



Longer. Faster. Forever.

Feature Selection

Rais Seminar

Vangjush Komini

KTH

April 1, 2020

Machine Learning in Nutshell

“Machine learning research is part of research on artificial intelligence, seeking to provide knowledge to computers through data, observations and interacting with the world. That acquired knowledge allows computers to correctly generalize to new settings.”

Yoshua Bengio

“Machine learning research is part of research on artificial intelligence, seeking to provide knowledge to computers through data, observations and interacting with the world. That acquired knowledge allows computers to correctly generalize to new settings.”

Yoshua Bengio

What is machine?

“Machine learning research is part of research on artificial intelligence, seeking to provide knowledge to computers through data, observations and interacting with the world. That acquired knowledge allows computers to correctly generalize to new settings.”

Yoshua Bengio

“Machine learning research is part of research on artificial intelligence, seeking to provide knowledge to computers through data, observations and interacting with the world. That acquired knowledge allows computers to correctly generalize to new settings.”

Yoshua Bengio

“Machine learning research is part of research on artificial intelligence, seeking to provide knowledge to computers through data, observations and interacting with the world. That acquired knowledge allows computers to correctly generalize to new settings.”

Yoshua Bengio

- Significant amounts of data are available or can be generated (either beforehand or dynamically)
- Other (analytical) solutions are too slow or infeasible
- Human expertise is absent or unexplainable
- Solutions change over time or need to be adapted
-

From DD3359-RL course KTH

- Significant amounts of data are available or can be generated (either beforehand or dynamically)
- Other (analytical) solutions are too slow or infeasible
 - Human expertise is absent or unexplainable
 - Solutions change over time or need to be adapted
 -

From DD3359-RL course KTH

- Significant amounts of data are available or can be generated (either beforehand or dynamically)
- Other (analytical) solutions are too slow or infeasible
- Human expertise is absent or unexplainable
- Solutions change over time or need to be adapted
-

From DD3359-RL course KTH

- Significant amounts of data are available or can be generated (either beforehand or dynamically)
- Other (analytical) solutions are too slow or infeasible
- Human expertise is absent or unexplainable
- Solutions change over time or need to be adapted

•

From DD3359-RL course KTH

- Significant amounts of data are available or can be generated (either beforehand or dynamically)
- Other (analytical) solutions are too slow or infeasible
- Human expertise is absent or unexplainable
- Solutions change over time or need to be adapted
-

From DD3359-RL course KTH

Deep learning allows computational models that are composed of multiple processing layers to learn representations of data, with multiple levels of abstraction.

Y.LeCun

Deep learning allows computational models, that are composed of multiple processing layers to learn representations of data, with multiple levels of abstraction.

Y.LeCun

Deep learning allows computational models, that are composed of multiple processing layers to learn representations of data, with multiple levels of abstraction.

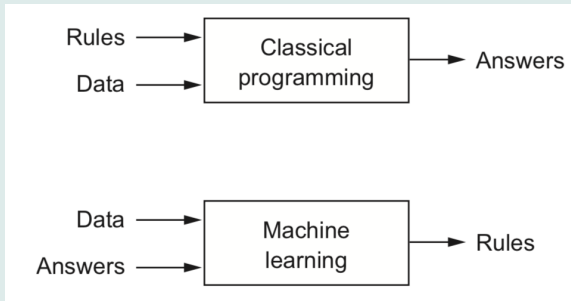
Y.LeCun

Deep learning allows computational models, that are composed of multiple processing layers to learn representations of data, with multiple levels of abstraction.

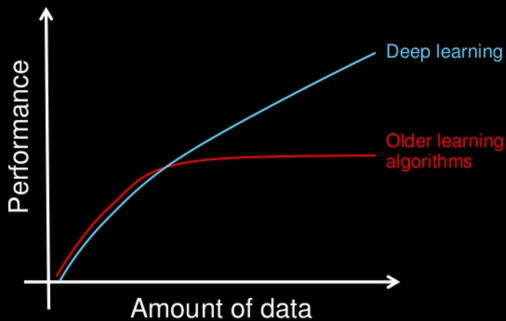
Y.LeCun

Deep learning allows computational models, that are composed of multiple processing layers to learn representations of data, with multiple levels of abstraction.

Y.LeCun



Why deep learning



How do data science techniques scale with amount of data?

Generalization vs Structuring

- Supervised learning tries to generalize over an massive amount of **structured** data.
- **Unsupervised** learning tries to **learn** the **structure** of a massive amount of data.
 - ▶ Clustering tries to bring together items with high similarity of **invariant** features.
 - ▶ Density estimation tries to model a probability distribution of the items influenced by the **invariant** features (Central Limit Theorem to be considered).
 - ▶ Dimensionality reduction find the a latent space where the **invariant** features prevail.
- **Semi-/Weakly- supervised** learning tries to **learn** the **scarcely** labeled data.
- Individual data is assumed to be **composed** of **core content** which is **invariant** from the **acquisition conditions** and the **non-core content dependent acquisition conditions**.

Generalization vs Structuring

- Supervised learning tries to generalize over an massive amount of **structured** data.
- **Unsupervised** learning tries to **learn** the **structure** of a massive amount of data.
 - ▶ **Clustering** tries to bring together items with high similarity of **invariant** features.
 - ▶ **Density estimation** tries to model a probability distribution of the items influenced by the **invariant** features (Central Limit Theorem to be considered).
 - ▶ **Dimensionality reduction** find the a latent space where the **invariant** features prevail.
- **Semi-/Weakly- supervised** learning tries to **learn** the **scarcely** labeled data.
- Individual data is assumed to be **composed** of **core content** which is **invariant** from the **acquisition conditions** and the **non-core content dependent acquisition conditions**.

Generalization vs Structuring

- Supervised learning tries to generalize over an massive amount of **structured** data.
- **Unsupervised** learning tries to **learn** the **structure** of a massive amount of data.
 - ▶ **Clustering** tries to bring together items with high similarity of **invariant** features.
 - ▶ **Density estimation** tries to model a probability distribution of the items influenced by the **invariant** features (Central Limit Theorem to be considered).
 - ▶ **Dimensionality reduction** find the a latent space where the **invariant** features prevail.
- **Semi-/Weakly- supervised** learning tries to **learn** the **scarcely** labeled data.
- Individual data is assumed to be **composed** of **core content** which is **invariant** from the **acquisition conditions** and the **non-core content dependent acquisition conditions**.

Generalization vs Structuring

- Supervised learning tries to generalize over an massive amount of **structured** data.
- **Unsupervised** learning tries to **learn** the **structure** of a massive amount of data.
 - ▶ **Clustering** tries to bring together items with high similarity of **invariant** features.
 - ▶ **Density estimation** tries to model a probability distribution of the items influenced by the **invariant** features (Central Limit Theorem to be considered).
 - ▶ **Dimensionality reduction** find the a latent space where the **invariant** features prevail.
- **Semi-/Weakly- supervised** learning tries to **learn** the **scarcely** labeled data.
- Individual data is assumed to be **composed** of **core content** which is **invariant** from the **acquisition conditions** and the **non-core content dependent acquisition conditions**.

Generalization vs Structuring

- Supervised learning tries to generalize over an massive amount of **structured** data.
- **Unsupervised** learning tries to **learn** the **structure** of a massive amount of data.
 - ▶ **Clustering** tries to bring together items with high similarity of **invariant** features.
 - ▶ **Density estimation** tries to model a probability distribution of the items influenced by the **invariant** features (Central Limit Theorem to be considered).
 - ▶ **Dimensionality reduction** find the a latent space where the **invariant** features prevail.
- **Semi-/Weakly- supervised** learning tries to **learn** the **scarcely** labeled data.
- Individual data is assumed to be **composed** of **core content** which is **invariant** from the **acquisition conditions** and the **non-core content dependent acquisition conditions**.

Generalization vs Structuring

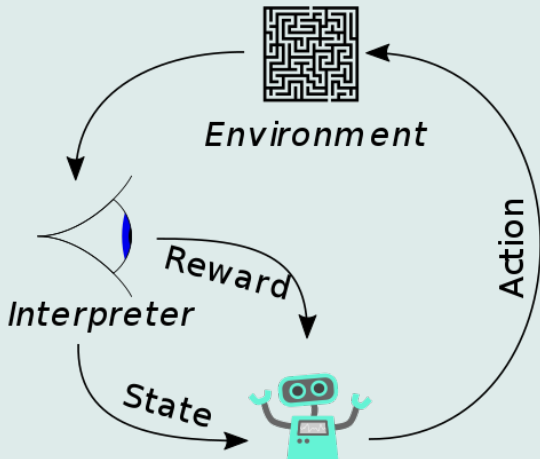
- Supervised learning tries to generalize over an massive amount of **structured** data.
- **Unsupervised** learning tries to **learn** the **structure** of a massive amount of data.
 - ▶ **Clustering** tries to bring together items with high similarity of **invariant** features.
 - ▶ **Density estimation** tries to model a probability distribution of the items influenced by the **invariant** features (Central Limit Theorem to be considered).
 - ▶ **Dimensionality reduction** find the a latent space where the **invariant** features prevail.
- **Semi-/Weakly- supervised** learning tries to **learn** the **scarcely** labeled data.
- Individual data is assumed to be **composed** of **core content** which is **invariant** from the **acquisition conditions** and the **non-core content dependent acquisition conditions**.

Generalization vs Structuring

- Supervised learning tries to generalize over an massive amount of **structured** data.
- **Unsupervised** learning tries to **learn** the **structure** of a massive amount of data.
 - ▶ **Clustering** tries to bring together items with high similarity of **invariant** features.
 - ▶ **Density estimation** tries to model a probability distribution of the items influenced by the **invariant** features (Central Limit Theorem to be considered).
 - ▶ **Dimensionality reduction** find the a latent space where the **invariant** features prevail.
- **Semi-/Weakly- supervised** learning tries to **learn** the **scarcely** labeled data.
- Individual data is assumed to be **composed** of **core content** which is **invariant** from the **acquisition conditions** and the **non-core** content **dependent acquisition conditions**.

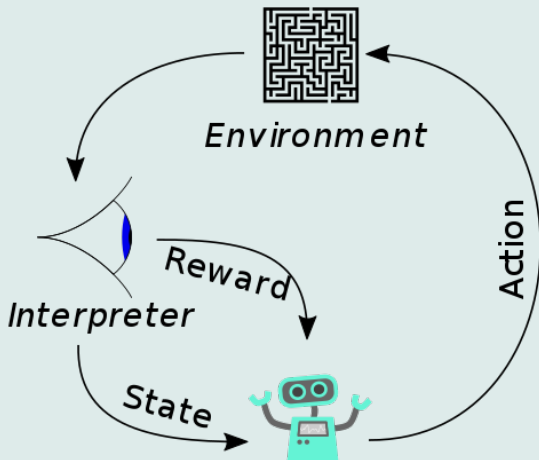
Other type

- Reinforcement learning: *Trying to generalize of a series of sequential observations from an environment by learning a policy generates a given action from the state, such as the cumulative rewards (sparse in time) are maximized.*
- More biologically plausible approach.



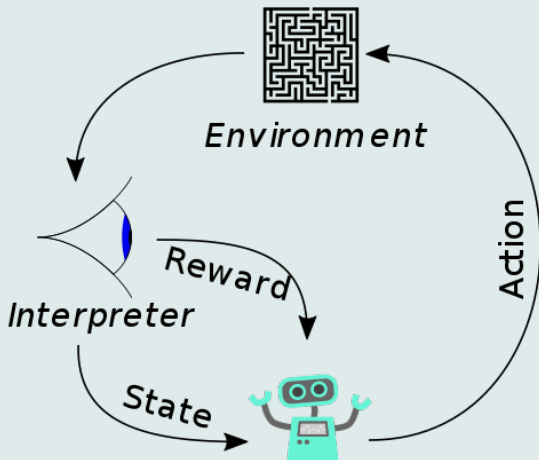
Other type

- Reinforcement learning: Trying to generalize of a series of sequential observations from an environment by learning a policy generates a given action from the state, such as the cumulative rewards (sparse in time) are maximized.
- More biologically plausible approach.



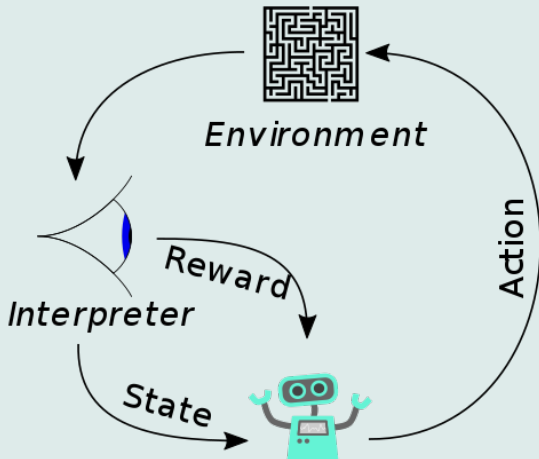
Other type

- Reinforcement learning: Trying to generalize of a series of sequential observations from an environment by learning a policy generates a given action from the state, such as the cumulative rewards (sparse in time) are maximized.
- More biologically plausible approach.



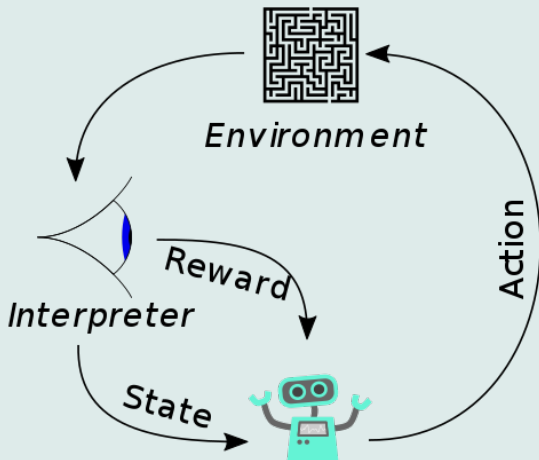
Other type

- Reinforcement learning: Trying to generalize of a series of sequential observations from an environment by learning a policy generates a given action from the state, such as the cumulative rewards (sparse in time) are maximized.
- More biologically plausible approach.



Other type

- Reinforcement learning: Trying to generalize of a series of sequential observations from an environment by learning a policy generates a given action from the state, such as the cumulative rewards (sparse in time) are maximized.
- More biologically plausible approach.



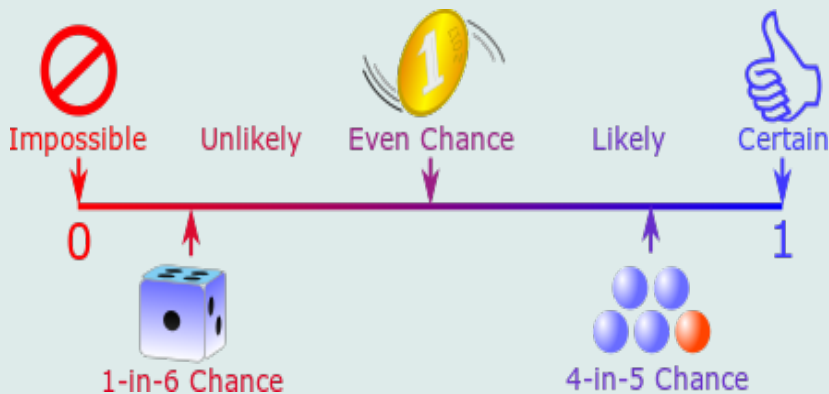
Treating the **concept** as a mathematical computable entity, and **sampling** a lot of data from the this entity, and use these empirical data as a proxy.

Treating the **concept** as a mathematical computable entity, and **sampling** a lot of data from the this entity, and use these empirical data as a proxy.

Treating the **concept** as a mathematical computable entity, and **sampling** a lot of data from the this entity, and use these empirical data as a proxy.

Key Mathematical Ingredients

- Probability: the calculus of **uncertainty** computation.
- Calculus: the science of **continuity** that is at continuous change.
- Algebra: the science of **multidimensional** hyper-space.
- Graphs: the science of **ontological** entities.
-



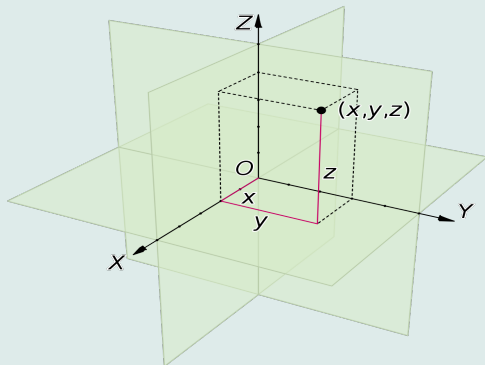
Key Mathematical Ingredients

- Probability: the calculus of **uncertainty** computation.
- Calculus: the science of **continuity** that is at continuous **change**.
- Algebra: the science of **multidimensional** hyper-space.
- Graphs: the science **ontological** entities.
-



Key Mathematical Ingredients

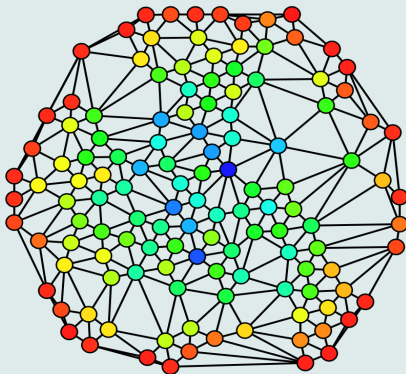
- Probability: the calculus of **uncertainty** computation.
- Calculus: the science of **continuity** that is at continuous **change**.
- Algebra: the science of **multidimensional** hyper-space.
- Graphs: the science **ontological** entities.
-



Key Mathematical Ingredients

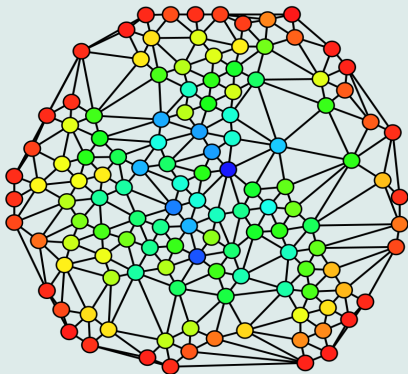
- Probability: the calculus of **uncertainty** computation.
- Calculus: the science of **continuity** that is at continuous **change**.
- Algebra: the science of **multidimensional** hyper-space.
- Graphs: the science **ontological** entities.

.....



Key Mathematical Ingredients

- Probability: the calculus of **uncertainty** computation.
- Calculus: the science of **continuity** that is at continuous **change**.
- Algebra: the science of **multidimensional** hyper-space.
- Graphs: the science **ontological** entities.
-

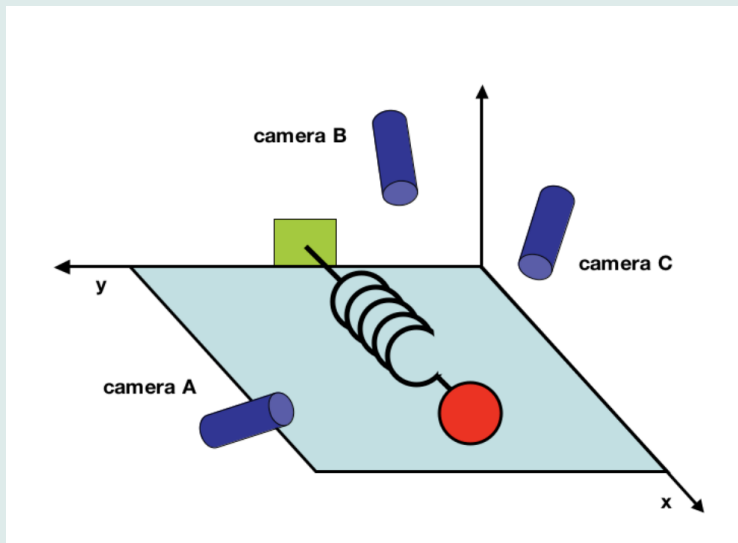


Feature selection



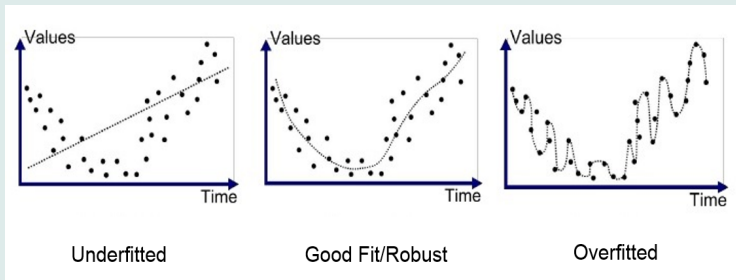
Which camera is better one

racefox



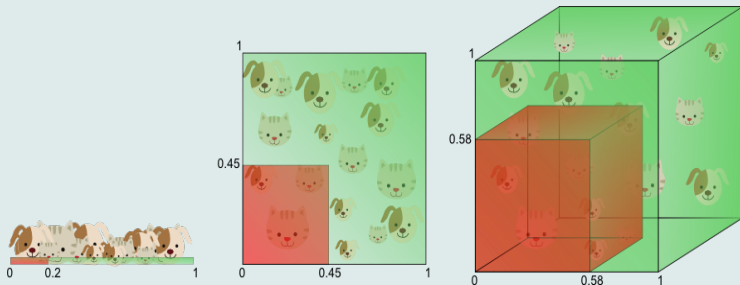
Why we need feature selection

- reducing overfitting
- overcoming the curse of dimensionality
- shorter training time
- improve the interpretability of the methods
- Other?



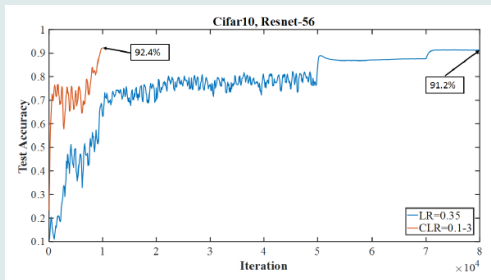
Why we need feature selection

- reducing overfitting
- overcoming the curse of dimensionality
- shorter training time
- improve the interpretability of the methods
- Other?



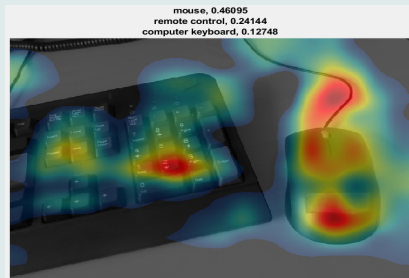
Why we need feature selection

- reducing overfitting
- overcoming the curse of dimensionality
- shorter training time
- improve the interpretability of the methods
- Other?



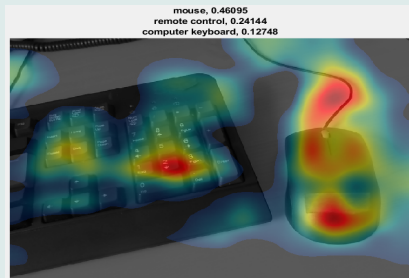
Why we need feature selection

- reducing overfitting
- overcoming the curse of dimensionality
- shorter training time
- improve the interpretability of the methods
- Other?



Why we need feature selection

- reducing overfitting
- overcoming the curse of dimensionality
- shorter training time
- improve the interpretability of the methods
- Other?



- Let $X \subset R^d$ be the domain of covariates.
- Let $Y \subset 0, 1$ be the domain of responses (labels).
- Given n i.i.d data pairs $\{(x_i, y_i), = 1, 2, \dots, d\}$, with unknown distribution $P(X, Y)$
- Select a subset of X that best predict Y .

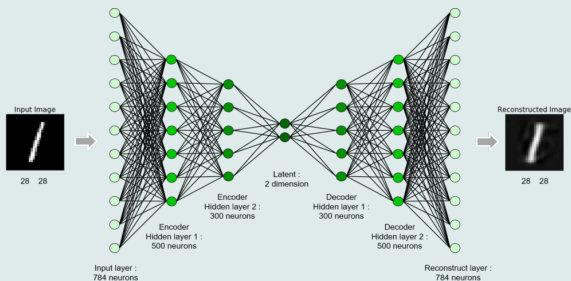
- Let $X \subset R^d$ be the domain of covariates.
- Let $Y \subset 0, 1$ be the domain of responses (labels).
- Given n i.i.d data pairs $\{(x_i, y_i), = 1, 2, \dots, d\}$, with unknown distribution $P(X, Y)$
- Select a subset of X that best predict Y .

- Let $X \subset R^d$ be the domain of covariates.
- Let $Y \subset 0, 1$ be the domain of responses (labels).
- Given n i.i.d data pairs $\{(x_i, y_i), = 1, 2, \dots, d\}$, with unknown distribution $P(X, Y)$
- Select a subset of X that best predict Y .

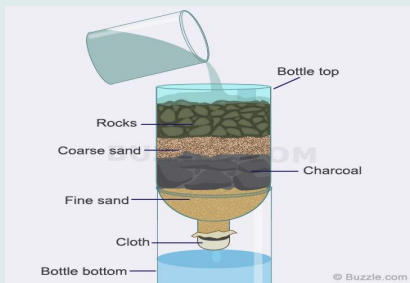
- Let $X \subset R^d$ be the domain of covariates.
- Let $Y \subset 0, 1$ be the domain of responses (labels).
- Given n i.i.d data pairs $\{(x_i, y_i), = 1, 2, \dots, d\}$, with unknown distribution $P(X, Y)$
- Select a subset of X that best predict Y .

Concrete autoencoder, Balin et al 2019

- Utilization of autoencoders, for distillation of predictive features.
- Latent space could be any type of mathematical entity.
- Reparameterization enables back-propagation in random variables.
- Concrete autoencoder is still an autoencoder.

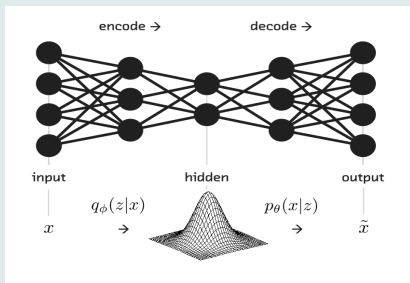


- Utilization of autoencoders, for distillation of predictive features.
- Latent space could be any type of mathematical entity.
- Reparameterization enables back-propagation in random variables.
- Concrete autoencoder is still an autoencoder.

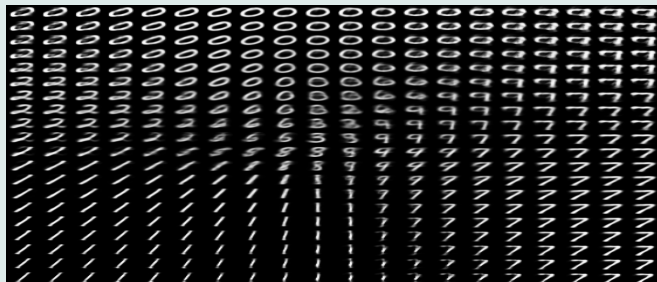


Concrete autoencoder

- Utilization of autoencoders, for distillation of predictive features.
- Latent space could be any type of mathematical entity.
- Reparameterization enables back-propagation in random variables.
- Concrete autoencoder is still an autoencoder.

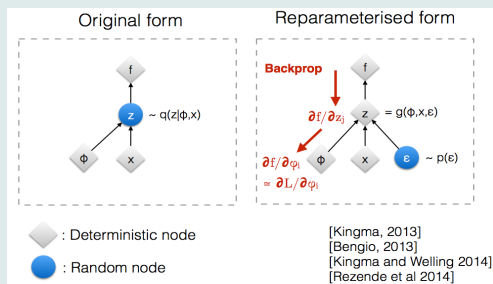


- Utilization of autoencoders, for distillation of predictive features.
- Latent space could be any type of mathematical entity.
- Reparameterization enables back-propagation in random variables.
- Concrete autoencoder is still an autoencoder.



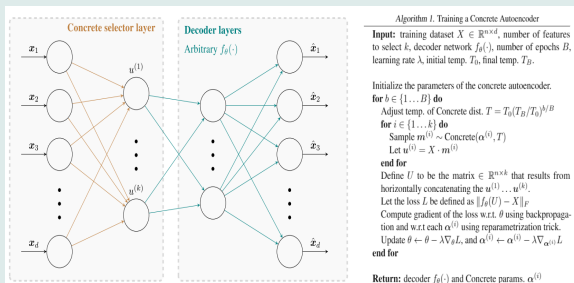
Concrete autoencoder

- Utilization of autoencoders, for distillation of predictive features.
- Latent space could be any type of mathematical entity.
- Reparameterization enables back-propagation in random variables.
- Concrete autoencoder is still an autoencoder.

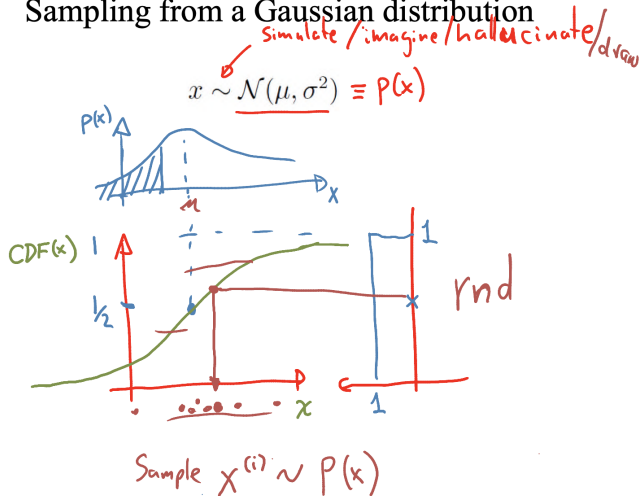


Concrete autoencoder

- Utilization of autoencoders, for distillation of predictive features.
- Latent space could be any type of mathematical entity.
- Reparameterization enables back-propagation in random variables.
- Concrete autoencoder is still an autoencoder.

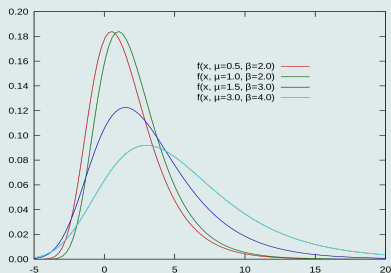


Sampling from a Gaussian distribution

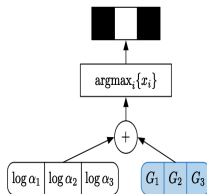
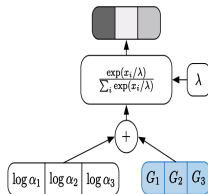


Slides taken from Nando.D.F course Machine Learning: 2014-2015

- Sample and array g_j from a Gumbel distribution $F(x; \mu, \beta) = e^{-e^{-\frac{x-\mu}{\beta}}}$
- Relaxation of the discrete variables $m_j = \frac{e^{(\log(\alpha_j) + g_j)/\lambda}}{\sum_{k=1}^d e^{(\log(\alpha_k) + g_k)/\lambda}}$
- One-hot encoding distribution $m_j = 1$, with probability $= \frac{\alpha_j}{\sum_{k=1}^d \alpha_k}$
- Temperature modulation $\lambda(\text{epoch}) = \lambda_{\text{initial}} \left\{ \frac{\lambda_{\text{final}}}{\lambda_{\text{initial}}} \right\}^{\left\{ \frac{\text{epoch}}{\text{TotalNrEpochs}} \right\}}$



- Sample and array g_j from a Gumbel distribution $F(x; \mu, \beta) = e^{-e^{-\frac{x-\mu}{\beta}}}$
- Relaxation of the discrete variables $m_j = \frac{e^{(\log(\alpha_j) + g_j)/\lambda}}{\sum_{k=1}^d e^{(\log(\alpha_k) + g_k)/\lambda}}$
- One-hot encoding distribution $m_j = 1$, with, probability $= \frac{\alpha_j}{\sum_{k=1}^d \alpha_k}$
- Temperature modulation $\lambda(\text{epoch}) = \lambda_{\text{initial}} \left\{ \frac{\lambda_{\text{final}}}{\lambda_{\text{initial}}} \right\}^{\left\{ \frac{\text{epoch}}{\text{TotalNrEpochs}} \right\}}$

(a) Discrete(α)(b) Concrete(α, λ)

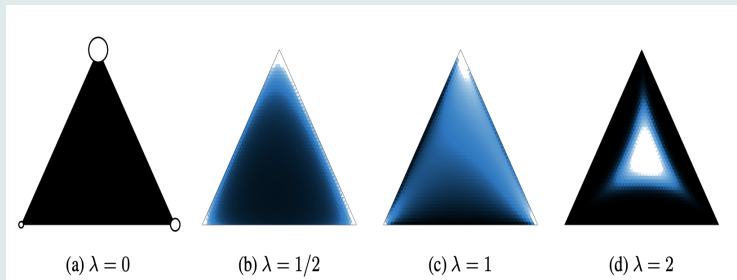
Concrete distribution, Maddison et al 2017

- Sample and array g_j from a Gumbel distribution $F(x; \mu, \beta) = e^{-e^{-\frac{x-\mu}{\beta}}}$
- Relaxation of the discrete variables $m_j = \frac{e^{(\log(\alpha_j) + g_j) / \lambda}}{\sum_{k=1}^d e^{(\log(\alpha_k) + g_k) / \lambda}}$
- One-hot encoding distribution $m_j = 1$, with, probability $= \frac{\alpha_j}{\sum_{k=1}^d \alpha_k}$
- Temperature modulation $\lambda(\text{epoch}) = \lambda_{\text{initial}} \left\{ \frac{\lambda_{\text{final}}}{\lambda_{\text{initial}}} \right\}^{\left\{ \frac{\text{epoch}}{\text{TotalNrEpochs}} \right\}}$

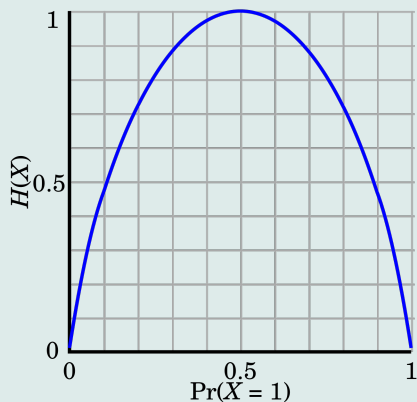
OneHot Encoding

workclass	State-gov	Self-emp-not-inc	Private
State-gov	1	0	0
Self-emp-not-inc	0	1	0
Private	0	0	1
Private	0	0	1
Private	0	0	1

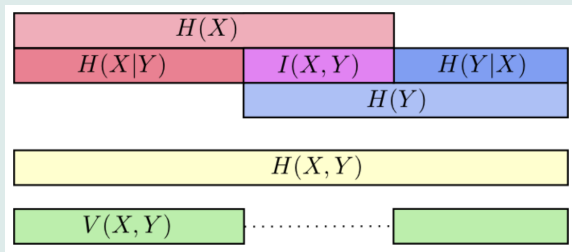
- Sample and array g_j from a Gumbel distribution $F(x; \mu, \beta) = e^{-e^{-\frac{x-\mu}{\beta}}}$
- Relaxation of the discrete variables $m_j = \frac{e^{(\log(\alpha_j)+g_j)/\lambda}}{\sum_{k=1}^d e^{(\log(\alpha_k)+g_k)/\lambda}}$
- One-hot encoding distribution $m_j = 1$, with, probability $= \frac{\alpha_j}{\sum_{k=1}^d \alpha_k}$
- Temperature modulation $\lambda(\text{epoch}) = \lambda_{\text{initial}} \left\{ \frac{\lambda_{\text{final}}}{\lambda_{\text{initial}}} \right\}^{\left\{ \frac{\text{epoch}}{\text{TotalNrEpochs}} \right\}}$



Variational Information Maximization for Feature Selection Gao et al

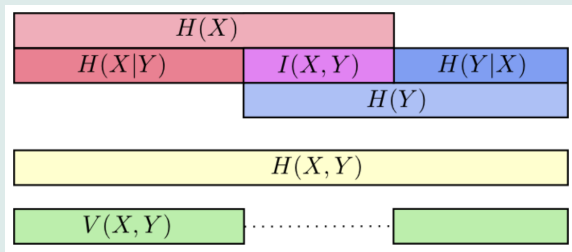


$$\text{Entropy : } H(X) = - \sum_i p_X(x_i) * \log(p_X(x_i)) \quad (1)$$



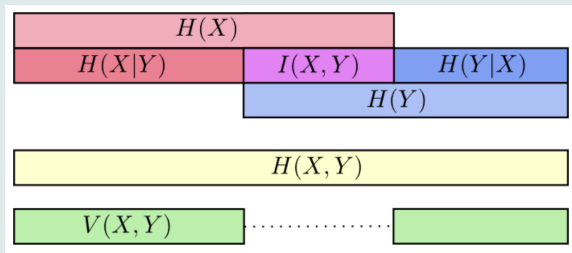
<https://colah.github.io/>

- $I(X) = H(X) + H(Y) - H(X, Y)$
- $V(X, Y) = H(X, Y) - I(X, Y)$
- $D_{KL}(P_X || Q_X) = \int_{-\infty}^{+\infty} p(x) \log \left\{ \frac{p(x)}{q(x)} \right\} dx$



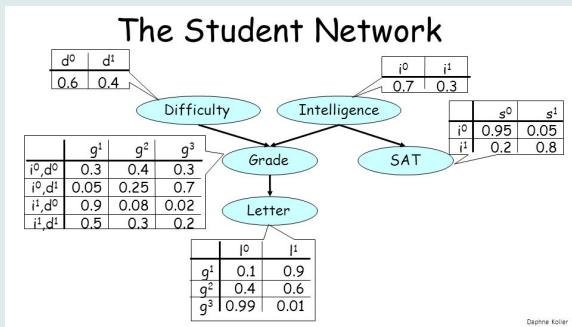
<https://colah.github.io/>

- $I(X) = H(X) + H(Y) - H(X, Y)$
- $VI(X, Y) = H(X, Y) - I(X, Y)$
- $D_{KL}(P_X || Q_X) = \int_{-\infty}^{+\infty} p(x) \log \left\{ \frac{p(x)}{q(x)} \right\} dx$

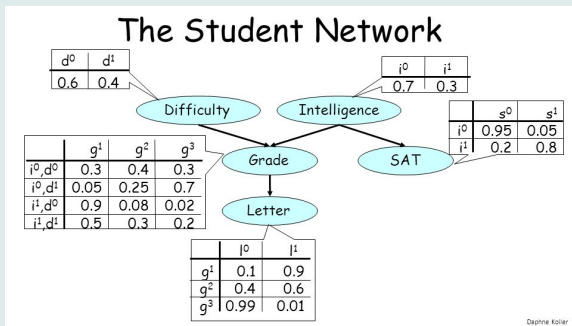


<https://colah.github.io/>

- $I(X) = H(X) + H(Y) - H(X, Y)$
- $VI(X, Y) = H(X, Y) - I(X, Y)$
- $D_{KL}(P_X || Q_X) = \int_{-\infty}^{+\infty} p(x) \log \left\{ \frac{p(x)}{q(x)} \right\} dx$



- Bayes : $P(\theta | \mathbf{Data}) = \frac{P(\theta, \mathbf{Data})}{P(\mathbf{Data})} = \frac{P(\theta) * P(\mathbf{Data} | \theta)}{P(\mathbf{Data})}$
- $P(X|Y) = P(X) \Rightarrow P(X, Y) = P(X)P(Y)$



- Bayes : $P(\theta | \text{Data}) = \frac{P(\theta, \text{Data})}{P(\text{Data})} = \frac{P(\theta) * P(\text{Data} | \theta)}{P(\text{Data})}$
- $P(X | Y) = P(X) \Rightarrow P(X, Y) = P(X)P(Y)$

- We would like a subset T of size (m)
s.t the remaining $S \setminus T$ are conditionally independent given T .
- This dependency is quantified by $Q : 2^d \rightarrow [0, \infty)$ such that:
 $Q(T)=0$ iff $X_{S \setminus T} \perp Y | X_T$
 $Q(T) \geq Q(S)$ whenever $T \subset S$

$$\min_{T:|T|=m} Q(T) \tag{2}$$

- We would like a subset T of size (m)
s.t the remaining $S \setminus T$ are conditionally independent given T .
- This dependency is quantified by $Q : 2^d \rightarrow [0, \infty)$ such that:

$$Q(T)=0 \text{ iff } X_{S \setminus T} \perp Y | X_T$$

$$Q(T) \geq Q(S) \text{ whenever } T \subset S$$

$$\min_{T:|T|=m} Q(T) \tag{2}$$

- We would like a subset T of size (m)
s.t the remaining $S \setminus T$ are conditionally independent given T .
- This dependency is quantified by $Q : 2^d \rightarrow [0, \infty)$ such that:
 $Q(T)=0$ iff $X_{S \setminus T} \perp Y | X_T$

$Q(T) \geq Q(S)$ whenever $T \subset S$

$$\min_{T:|T|=m} Q(T) \tag{2}$$

- We would like a subset T of size (m)
s.t the remaining $S \setminus T$ are conditionally independent given T .
- This dependency is quantified by $Q : 2^d \rightarrow [0, \infty)$ such that:
 $Q(T)=0$ iff $X_{S \setminus T} \perp Y | X_T$
 $Q(T) \geq Q(S)$ whenever $T \subset S$

$$\min_{T:|T|=m} Q(T) \tag{2}$$

Feature just random variables

- $T = \operatorname{argmax}_T \left\{ I\{(x_1, x_2, \dots, x_T), y\} \right\}$ NP-hard direct solution
- Forward Feature Selection : $t_{\text{step}=t} = \operatorname{argmax}_{i \in S^{t-1}} \left\{ I(x_{S^{t-1} \cup i} : y) \right\}$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y | x_{S^{t-1}})$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y) - I(x_i : x_{S^{t-1}}) + I(x_i : x_{S^{t-1}} | y)$
- $= I(x_{S^{t-1}} : y) + I(x_i : y) - (H(x_{S^{t-1}}) - H(x_{S^{t-1}} | x_i)) + (H(x_{S^{t-1}} | y) - H(x_{S^{t-1}} | x_i, y))$
- $t = \operatorname{argmax}_{i \in S^{t-1}} \left\{ I(x_i : y) + H(x_{S^{t-1}} | x_i) - H(x_{S^{t-1}} | x_i, y) \right\}$
- $H(x_{S^{t-1}} | x_i) \approx \sum_{k=1}^{t-1} H(x_k | x_i)$
- $H(x_{S^{t-1}} | x_i, y) \approx \sum_{k=1}^{t-1} H(x_k | x_i, y)$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $T = \operatorname{argmax}_T \left\{ I\{(x_1, x_2, \dots, x_T), y\} \right\}$ NP-hard direct solution
- **Forward Feature Selection** : $t_{\text{step}=t} = \operatorname{argmax}_{i \in S^{t-1}} \left\{ I(x_{S^{t-1} \cup i} : y) \right\}$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y | x_{S^{t-1}})$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y) - I(x_i : x_{S^{t-1}}) + I(x_i : x_{S^{t-1}} | y)$
- $= I(x_{S^{t-1}} : y) + I(x_i : y) - (H(x_{S^{t-1}}) - H(x_{S^{t-1}} | x_i)) + (H(x_{S^{t-1}} | y) - H(x_{S^{t-1}} | x_i, y))$
- $t = \operatorname{argmax}_{i \in S^{t-1}} \left\{ I(x_i : y) + H(x_{S^{t-1}} | x_i) - H(x_{S^{t-1}} | x_i, y) \right\}$
- $H(x_{S^{t-1}} | x_i) \approx \sum_{k=1}^{t-1} H(x_k | x_i)$
- $H(x_{S^{t-1}} | x_i, y) \approx \sum_{k=1}^{t-1} H(x_k | x_i, y)$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $T = \operatorname{argmax}_T \left\{ I\{(x_1, x_2, \dots, x_T), y\} \right\}$ NP-hard direct solution
- **Forward Feature Selection** : $t_{\text{step}=t} = \operatorname{argmax}_{i \notin S^{t-1}} \left\{ I(x_{S^{t-1} \cup i} : y) \right\}$
 - $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y | x_{S^{t-1}})$
 - $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y) - I(x_i : x_{S^{t-1}}) + I(x_i : x_{S^{t-1}} | y)$
 - $= I(x_{S^{t-1}} : y) + I(x_i : y) - (H(x_{S^{t-1}}) - H(x_{S^{t-1}} | x_i)) + (H(x_{S^{t-1}} | y) - H(x_{S^{t-1}} | x_i, y))$
 - $t = \operatorname{argmax}_{i \notin S^{t-1}} \left\{ I(x_i : y) + H(x_{S^{t-1}} | x_i) - H(x_{S^{t-1}} | x_i, y) \right\}$
 - $H(x_{S^{t-1}} | x_i) \approx \sum_{k=1}^{t-1} H(x_k | x_i)$
 - $H(x_{S^{t-1}} | x_i, y) \approx \sum_{k=1}^{t-1} H(x_k | x_i, y)$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $T = \operatorname{argmax}_T \left\{ I\{(x_1, x_2, \dots, x_T), y\} \right\}$ NP-hard direct solution
- **Forward Feature Selection** : $t_{\text{step}=t} = \operatorname{argmax}_{i \notin S^{t-1}} \left\{ I(x_{S^{t-1} \cup i} : y) \right\}$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y | x_{S^{t-1}})$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y) - I(x_i : x_{S^{t-1}}) + I(x_i : x_{S^{t-1}} | y)$
- $= I(x_{S^{t-1}} : y) + I(x_i : y) - (H(x_{S^{t-1}}) - H(x_{S^{t-1}} | x_i)) + (H(x_{S^{t-1}} | y) - H(x_{S^{t-1}} | x_i, y))$
- $t = \operatorname{argmax}_{i \notin S^{t-1}} \left\{ I(x_i : y) + H(x_{S^{t-1}} | x_i) - H(x_{S^{t-1}} | x_i, y) \right\}$
- $H(x_{S^{t-1}} | x_i) \approx \sum_{k=1}^{t-1} H(x_k | x_i)$
- $H(x_{S^{t-1}} | x_i, y) \approx \sum_{k=1}^{t-1} H(x_k | x_i, y)$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $T = \operatorname{argmax}_T \left\{ I\{(x_1, x_2, \dots, x_T), y\} \right\}$ NP-hard direct solution
- **Forward Feature Selection** : $t_{\text{step}=t} = \operatorname{argmax}_{i \in S^{t-1}} \left\{ I(x_{S^{t-1} \cup i} : y) \right\}$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y | x_{S^{t-1}})$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y) - I(x_i : x_{S^{t-1}}) + I(x_i : x_{S^{t-1}} | y)$
- $= I(x_{S^{t-1}} : y) + I(x_i : y) - (H(x_{S^{t-1}}) - H(x_{S^{t-1}} | x_i)) + (H(x_{S^{t-1}} | y) - H(x_{S^{t-1}} | x_i, y))$
- $t = \operatorname{argmax}_{i \in S^{t-1}} \left\{ I(x_i : y) + H(x_{S^{t-1}} | x_i) - H(x_{S^{t-1}} | x_i, y) \right\}$
- $H(x_{S^{t-1}} | x_i) \approx \sum_{k=1}^{t-1} H(x_k | x_i)$
- $H(x_{S^{t-1}} | x_i, y) \approx \sum_{k=1}^{t-1} H(x_k | x_i, y)$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $T = \operatorname{argmax}_T \left\{ I\{(x_1, x_2, \dots, x_T), y\} \right\}$ NP-hard direct solution
- **Forward Feature Selection** : $t_{\text{step}=t} = \operatorname{argmax}_{i \in S^{t-1}} \left\{ I(x_{S^{t-1} \cup i} : y) \right\}$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y | x_{S^{t-1}})$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y) - I(x_i : x_{S^{t-1}}) + I(x_i : x_{S^{t-1}} | y)$
- $= I(x_{S^{t-1}} : y) + I(x_i : y) - (H(x_{S^{t-1}}) - H(x_{S^{t-1}} | x_i)) + (H(x_{S^{t-1}} | y) - H(x_{S^{t-1}} | x_i, y))$
- $t = \operatorname{argmax}_{i \in S^{t-1}} \left\{ I(x_i : y) + H(x_{S^{t-1}} | x_i) - H(x_{S^{t-1}} | x_i, y) \right\}$
- $H(x_{S^{t-1}} | x_i) \approx \sum_{k=1}^{t-1} H(x_k | x_i)$
- $H(x_{S^{t-1}} | x_i, y) \approx \sum_{k=1}^{t-1} H(x_k | x_i, y)$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $T = \operatorname{argmax}_T \left\{ I\{(x_1, x_2, \dots, x_T), y\} \right\}$ NP-hard direct solution
- **Forward Feature Selection** : $t_{\text{step}=t} = \operatorname{argmax}_{i \notin S^{t-1}} \left\{ I(x_{S^{t-1} \cup i} : y) \right\}$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y | x_{S^{t-1}})$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y) - I(x_i : x_{S^{t-1}}) + I(x_i : x_{S^{t-1}} | y)$
- $= I(x_{S^{t-1}} : y) + I(x_i : y) - (H(x_{S^{t-1}}) - H(x_{S^{t-1}} | x_i)) + (H(x_{S^{t-1}} | y) - H(x_{S^{t-1}} | x_i, y))$
- $t = \operatorname{argmax}_{i \notin S^{t-1}} \left\{ I(x_i : y) + H(x_{S^{t-1}} | x_i) - H(x_{S^{t-1}} | x_i, y) \right\}$
- $H(x_{S^{t-1}} | x_i) \approx \sum_{k=1}^{t-1} H(x_k | x_i)$
- $H(x_{S^{t-1}} | x_i, y) \approx \sum_{k=1}^{t-1} H(x_k | x_i, y)$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $T = \operatorname{argmax}_T \left\{ I \left\{ (x_1, x_2, \dots, x_T), y \right\} \right\}$ NP-hard direct solution
- **Forward Feature Selection** : $t_{\text{step}=t} = \operatorname{argmax}_{i \notin S^{t-1}} \left\{ I(x_{S^{t-1} \cup i} : y) \right\}$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y | x_{S^{t-1}})$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y) - I(x_i : x_{S^{t-1}}) + I(x_i : x_{S^{t-1}} | y)$
- $= I(x_{S^{t-1}} : y) + I(x_i : y) - (H(x_{S^{t-1}}) - H(x_{S^{t-1}} | x_i)) + (H(x_{S^{t-1}} | y) - H(x_{S^{t-1}} | x_i, y))$
- $t = \operatorname{argmax}_{i \notin S^{t-1}} \left\{ I(x_i : y) + H(x_{S^{t-1}} | x_i) - H(x_{S^{t-1}} | x_i, y) \right\}$
- $H(x_{S^{t-1}} | x_i) \approx \sum_{k=1}^{t-1} H(x_k | x_i)$
- $H(x_{S^{t-1}} | x_i, y) \approx \sum_{k=1}^{t-1} H(x_k | x_i, y)$

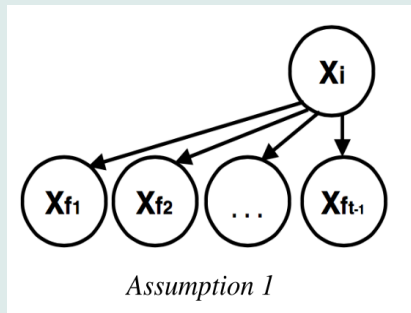
Not intended to be understood at a single slide. Check reference for further understanding.

- $T = \operatorname{argmax}_T \left\{ I \left\{ (x_1, x_2, \dots, x_T), y \right\} \right\}$ NP-hard direct solution
- **Forward Feature Selection** : $t_{\text{step}=t} = \operatorname{argmax}_{i \notin S^{t-1}} \left\{ I(x_{S^{t-1} \cup i} : y) \right\}$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y | x_{S^{t-1}})$
- $I(x_{S^{t-1} \cup i} : y) = I(x_{S^{t-1}} : y) + I(x_i : y) - I(x_i : x_{S^{t-1}}) + I(x_i : x_{S^{t-1}} | y)$
- $= I(x_{S^{t-1}} : y) + I(x_i : y) - (H(x_{S^{t-1}}) - H(x_{S^{t-1}} | x_i)) + (H(x_{S^{t-1}} | y) - H(x_{S^{t-1}} | x_i, y))$
- $t = \operatorname{argmax}_{i \notin S^{t-1}} \left\{ I(x_i : y) + H(x_{S^{t-1}} | x_i) - H(x_{S^{t-1}} | x_i, y) \right\}$
- $H(x_{S^{t-1}} | x_i) \approx \sum_{k=1}^{t-1} H(x_k | x_i)$
- $H(x_{S^{t-1}} | x_i, y) \approx \sum_{k=1}^{t-1} H(x_k | x_i, y)$

Not intended to be understood at a single slide. Check reference for further understanding.

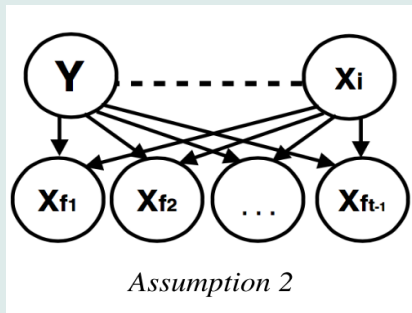
Feature just random variables

- Assumption1: Feature Independent : $P(x_{S_{t-1}} | x_i) = \prod_{k=1}^{t-1} P(x_k | x_i)$
- Assumption2: Class-Conditioned Independent
 $P(x_{S_{t-1}} | x_i, y) = \prod_{k=1}^{t-1} P(x_k | x_i, y)$
- These have only one common structure fulfillment.
- Contradiction when both are met: $I(x_i; Y) > I(x_1, x_2, \dots, x_{t-1}; Y)$



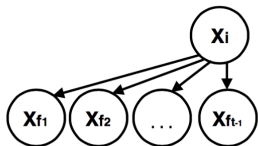
Feature just random variables

- Assumption1: Feature Independent : $P(x_{t-1}|x) = \prod_{k=1}^{t-1} P(x_k|x)$
- Assumption2: Class-Conditioned Independent
: $P(x_{t-1}|x_i, y) = \prod_{k=1}^{t-1} P(x_k|x_i, y)$
- These have only one common structure fulfillment.
- Contradiction when both are met: $I(x_i; Y) > I(x_1, x_2, \dots, x_{t-1}; Y)$

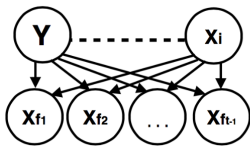


Feature just random variables

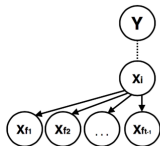
- Assumption1: Feature Independent : $P(x_{q-1}|x) = \prod_{i=1}^{q-1} P(x_i|x)$
- Assumption2: Class-Conditioned Independent
: $P(x_{q-1}|x, y) = \prod_{i=1}^{q-1} P(x_i|x, y)$
- These have only one common structure fulfillment.
- Contradiction when both are met: $I(x_i : y) > I(x_1, x_2, \dots, x_{i-1} : y)$



Assumption 1



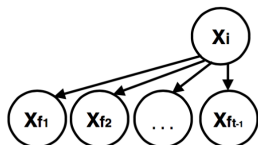
Assumption 2



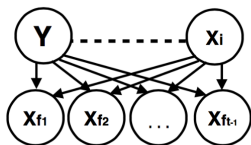
Satisfying both Assumption 1 and Assumption 2

Feature just random variables

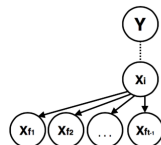
- Assumption1: Feature Independent : $P(x_{1:t-1}|y) = \prod_{i=1}^{t-1} P(x_i|y)$
- Assumption2: Class-Conditioned Independent
: $P(x_{1:t-1}|y, x) = \prod_{i=1}^{t-1} P(x_i|x, y)$
- These have only one common structure fulfillment.
- Contradiction when both are met: $I(x_i : y) > I(x_1, x_2, \dots, x_{t-1} : y)$



Assumption 1



Assumption 2



Satisfying both Assumption 1 and Assumption 2

Feature just random variables

- $I(x, y) \geq H(x) + \langle \ln[q(x|y)] \rangle_{p(x,y)}$
- $S = \operatorname{argmax}_S \{ H(x_S) + \langle \ln[q(x_S|y)] \rangle_{p(x_S,y)} \}$
- Swap x with y: $I(x, y) \geq H(y) + \langle \ln[q(y|x)] \rangle_{p(x,y)} = \langle \ln \left\{ \frac{q(y|x)}{p(y)} \right\} \rangle_{p(x,y)}$
- $S^* = \operatorname{argmax}_S \left\{ \langle \ln \left\{ \frac{q(y|x_S)}{p(y)} \right\} \rangle_{p(x_S,y)} \right\}$
- $q(y|x_S) = \frac{q(x_S,y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{\sum_j q(x_S,y) \ln(y^j)}$
- $I(x_S, y) \geq \left\langle \ln \left\{ \frac{q(y|x_S)}{q(x_S)} \right\} \right\rangle_{p(x_S,y)} = I_{LB}(x_S : y)$
- $I(x_S, y) - I_{LB}(x_S : y) = \langle KL(p(y|x_S) || q(y|x_S)) \rangle_{p(x_S)}$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $I(x, y) \geq H(x) + \langle \ln[q(x|y)] \rangle_{p(x,y)}$
- $S = \operatorname{argmax}_S \{H(x_S) + \langle \ln(q(x_S|y)) \rangle_{p(x_S,y)}\}$
- Swap x with y: $I(x, y) \geq H(y) + \langle \ln[q(y|x)] \rangle_{p(x,y)} = \langle \ln\left\{\frac{q(y|x)}{p(y)}\right\} \rangle_{p(x,y)}$
- $S^* = \operatorname{argmax}_S \left\{ \langle \ln\left\{\frac{q(y|x_S)}{p(y)}\right\} \rangle_{p(x_S,y)} \right\}$
- $q(y|x_S) = \frac{q(x_S,y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{\sum_y q(x_S,y) \ln(y^*)}$
- $I(x_S, y) \geq \left\langle \ln\left\{\frac{q(y|x_S)}{q(x_S)}\right\} \right\rangle_{p(x_S,y)} = I_{LB}(x_S : y)$
- $I(x_S, y) - I_{LB}(x_S : y) = \langle KL(p(y|x_S) || q(y|x_S)) \rangle_{p(x_S)}$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $I(x, y) \geq H(x) + \langle \ln[q(x|y)] \rangle_{p(x,y)}$
- $S = \operatorname{argmax}_S \{ H(x_S) + \langle \ln[q(x_S|y)] \rangle_{p(x_S,y)} \}$
- Swap x with y: $I(x, y) \geq H(y) + \langle \ln[q(y|x)] \rangle_{p(x,y)} = \langle \ln \left\{ \frac{q(y|x)}{p(y)} \right\} \rangle_{p(x,y)}$
- $S^* = \operatorname{argmax}_S \left\{ \langle \ln \left\{ \frac{q(y|x_S)}{p(y)} \right\} \rangle_{p(x_S,y)} \right\}$
- $q(y|x_S) = \frac{q(x_S,y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{\sum_y q(x_S,y)p(y)}$
- $I(x_S, y) \geq \left\langle \ln \left\{ \frac{q(y|x_S)}{q(x_S)} \right\} \right\rangle_{p(x_S,y)} = I_{LB}(x_S : y)$
- $I(x_S, y) - I_{LB}(x_S : y) = \langle KL(p(y|x_S) || q(y|x_S)) \rangle_{p(x_S)}$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $I(x, y) \geq H(x) + \langle \ln[q(x|y)] \rangle_{p(x,y)}$
- $S = \operatorname{argmax}_S \{ H(x_S) + \langle \ln[q(x_S|y)] \rangle_{p(x_S,y)} \}$
- Swap x with y: $I(x, y) \geq H(y) + \langle \ln[q(y|x)] \rangle_{p(x,y)} = \langle \ln \left\{ \frac{q(y|x)}{p(y)} \right\} \rangle_{p(x,y)}$
- $S^* = \operatorname{argmax}_S \left\{ \langle \ln \left\{ \frac{q(y|x_S)}{p(y)} \right\} \rangle_{p(x_S,y)} \right\}$
- $q(y|x_S) = \frac{q(x_S,y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{\sum_x q(x_S,y) \ln(x^*)}$
- $I(x_S, y) \geq \left\langle \ln \left\{ \frac{q(y|x_S)}{q(x_S)} \right\} \right\rangle_{p(x_S,y)} = I_{LB}(x_S : y)$
- $I(x_S, y) - I_{LB}(x_S : y) = \langle KL(p(y|x_S) || q(y|x_S)) \rangle_{p(x_S)}$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $I(x, y) \geq H(x) + \langle \ln[q(x|y)] \rangle_{p(x,y)}$
- $S = \operatorname{argmax}_S \{ H(x_S) + \langle \ln[q(x_S|y)] \rangle_{p(x_S,y)} \}$
- Swap x with y: $I(x, y) \geq H(y) + \langle \ln[q(y|x)] \rangle_{p(x,y)} = \langle \ln \left\{ \frac{q(y|x)}{p(y)} \right\} \rangle_{p(x,y)}$
- $S^* = \operatorname{argmax}_S \left\{ \langle \ln \left\{ \frac{q(y|x_S)}{p(y)} \right\} \rangle_{p(x_S,y)} \right\}$
- $q(y|x_S) = \frac{q(x_S,y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{\sum_{y'} q(x_S|y')p(y')}$
- $I(x_S, y) \geq \left\langle \ln \left\{ \frac{q(y|x_S)}{q(x_S)} \right\} \right\rangle_{p(x_S,y)} = I_{LB}(x_S : y)$
- $I(x_S, y) - I_{LB}(x_S : y) = \langle KL(p(y|x_S) || q(y|x_S)) \rangle_{p(x_S)}$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $I(x, y) \geq H(x) + \langle \ln[q(x|y)] \rangle_{p(x,y)}$
- $S = \operatorname{argmax}_S \{ H(x_S) + \langle \ln[q(x_S|y)] \rangle_{p(x_S,y)} \}$
- Swap x with y: $I(x, y) \geq H(y) + \langle \ln[q(y|x)] \rangle_{p(x,y)} = \langle \ln \left\{ \frac{q(y|x)}{p(y)} \right\} \rangle_{p(x,y)}$
- $S^* = \operatorname{argmax}_S \left\{ \langle \ln \left\{ \frac{q(y|x_S)}{p(y)} \right\} \rangle_{p(x_S,y)} \right\}$
- $q(y|x_S) = \frac{q(x_S,y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{\sum_{y'} q(x_S|y')p(y')}$
- $I(x_S, y) \geq \left\langle \ln \left\{ \frac{q(y|x_S)}{q(x_S)} \right\} \right\rangle_{p(x_S,y)} = I_{LB}(x_S : y)$
- $I(x_S, y) - I_{LB}(x_S : y) = \langle KL(p(y|x_S) || q(y|x_S)) \rangle_{p(x_S)}$

Not intended to be understood at a single slide. Check reference for further understanding.

Feature just random variables

- $I(x, y) \geq H(x) + \langle \ln[q(x|y)] \rangle_{p(x,y)}$
- $S = \operatorname{argmax}_S \{ H(x_S) + \langle \ln[q(x_S|y)] \rangle_{p(x_S,y)} \}$
- Swap x with y: $I(x, y) \geq H(y) + \langle \ln[q(y|x)] \rangle_{p(x,y)} = \langle \ln \left\{ \frac{q(y|x)}{p(y)} \right\} \rangle_{p(x,y)}$
- $S^* = \operatorname{argmax}_S \left\{ \langle \ln \left\{ \frac{q(y|x_S)}{p(y)} \right\} \rangle_{p(x_S,y)} \right\}$
- $q(y|x_S) = \frac{q(x_S,y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{q(x_S)} = \frac{q(x_S,y)p(y)}{\sum_{y'} q(x_S|y')p(y')}$
- $I(x_S, y) \geq \left\langle \ln \left\{ \frac{q(y|x_S)}{p(y)} \right\} \right\rangle_{p(x_S,y)} = I_{LB}(x_S : y)$
- $I(x_S, y) - I_{LB}(x_S : y) = \langle KL(p(y|x_S) || q(y|x_S)) \rangle_{p(x_S)}$

Not intended to be understood at a single slide. Check reference for further understanding.

Variational feature selection under auto-regressive decomposition

- $q(x_s|y) = q(x_1|y) \prod_{t=2}^T q(x_t|x_{T \geq t}, y)$
- $I_{LB}(x_s : y) = \frac{1}{N} \sum_{x^k, y^k} \ln \frac{q(x_s^k|y^k)}{q(x_s^k)}$

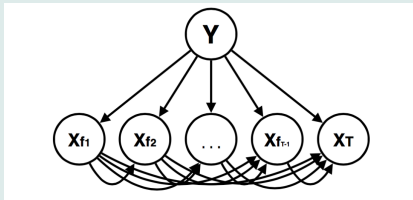


Figure 2: Auto-regressive decomposition for $q(x_S|y)$

- MI assesses the informativeness of features
- It requires a lot of observation if the dimensionality of the data is very high

Variational feature selection under auto-regressive decomposition

- $q(x_s|y) = q(x_1|y) \prod_{t=2}^T q(x_t|x_{T \geq t}, y)$
- $I_{LB}(x_s : y) = \frac{1}{N} \sum_{x^k, y^k} \ln \frac{\hat{q}(x_s^k|y^k)}{\hat{q}(x_s^{(k)})}$

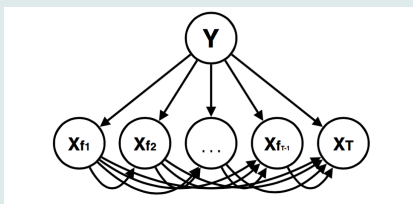


Figure 2: Auto-regressive decomposition for $q(x_s|y)$

- MI assesses the informativeness of features
- It requires a lot of observation if the dimensionality of the data is very high

Variational feature selection under auto-regressive decomposition

- $q(x_s|y) = q(x_1|y) \prod_{t=2}^T q(x_t|x_{T \geq t}, y)$
- $I_{LB}(x_s : y) = \frac{1}{N} \sum_{x^k, y^k} \ln \frac{\hat{q}(x_s^k|y^k)}{\hat{q}(x_s^{(k)})}$

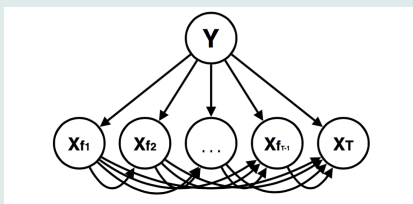


Figure 2: Auto-regressive decomposition for $q(x_s|y)$

- **MI assesses the informativeness of features**
- It requires a lot of observation if the dimensionality of the data is very high

Variational feature selection under auto-regressive decomposition

- $q(x_s|y) = q(x_1|y) \prod_{t=2}^T q(x_t|x_{T \geq t}, y)$
- $I_{LB}(x_s : y) = \frac{1}{N} \sum_{x^k, y^k} \ln \frac{\hat{q}(x_s^k|y^k)}{\hat{q}(x_s^{(k)})}$

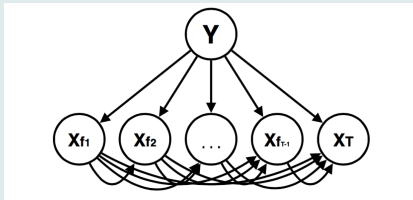


Figure 2: Auto-regressive decomposition for $q(x_s|y)$

- MI assesses the informativeness of features
- It requires a lot of observation if the dimensionality of the data is very high

Kernel Feature Selection via Conditional Covariance Minimization

Chen et al

- F is a class of functions from X to Y .
- L is a **loss function** defined by the user (MSE).
- Prediction error: $\epsilon_F = \inf_{f \in F} E_{X,Y} L(Y, f(X))$
- Solve the problem:

$$\min_{T:|T| \leq m} \epsilon_F(X_T) = \min_{T:|T| \leq m} \inf_{f \in F} E_{X,Y} \left\{ L(T, f(X)) \right\}$$

- F is a class of functions from X to Y .
- L is a **loss function** defined by the user (MSE).

- Prediction error: $\epsilon_F = \inf_{f \in F} E_{X,Y} L(Y, f(X))$

- Solve the problem:

$$\min_{T: |T| \leq m} \epsilon_F(X_T) = \min_{T: |T| \leq m} \inf_{f \in F} E_{X,Y} \left\{ L(T, f(X)) \right\}$$

- F is a class of functions from X to Y .
- L is a **loss function** defined by the user (MSE).
- Prediction error: $\epsilon_F = \inf_{f \in F} E_{X,Y} L(Y, f(X))$
- Solve the problem:

$$\min_{T:|T|\leq m} \epsilon_F(X_T) = \min_{T:|T|\leq m} \inf_{f \in F} E_{X,Y} \left\{ L(T, f(X)) \right\}$$

- F is a class of functions from X to Y .
- L is a **loss function** defined by the user (MSE).
- Prediction error: $\epsilon_F = \inf_{f \in F} E_{X,Y} L(Y, f(X))$
- Solve the problem:

$$\min_{T:|T|\leq m} \epsilon_F(X_T) = \min_{T:|T|\leq m} \inf_{f \in F} E_{X,Y} \left\{ L(T, f(X)) \right\}$$

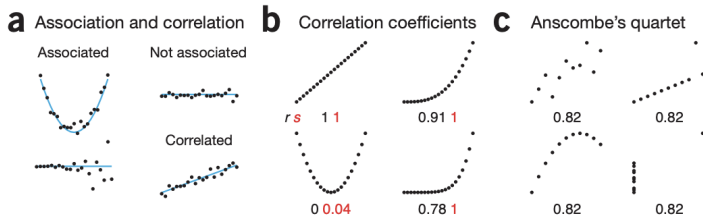
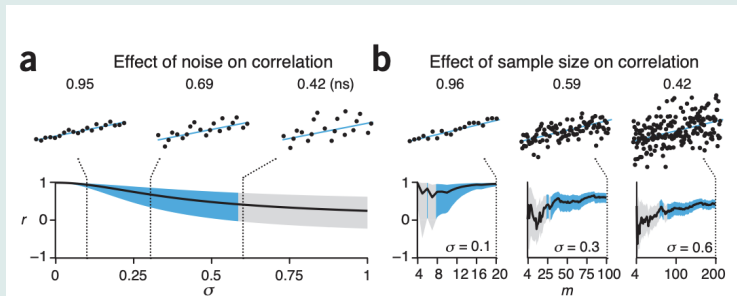


Figure 1 | Correlation is a type of association and measures increasing or decreasing trends quantified using correlation coefficients. **(a)** Scatter plots of associated (but not correlated), non-associated and correlated variables. In the lower association example, variance in y is increasing with x . **(b)** The Pearson correlation coefficient (r , black) measures linear trends, and the Spearman correlation coefficient (s , red) measures increasing or decreasing trends. **(c)** Very different data sets may have similar r values. Descriptors such as curvature or the presence of outliers can be more specific.

Altman et al

- Correlation: How much variance is explained $\rightarrow r_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y}$
- Covariance: X and Y co-vary $\rightarrow \text{cov}(X,Y) = r_{X,Y} \cdot \sigma_X \cdot \sigma_Y$



Altman et al

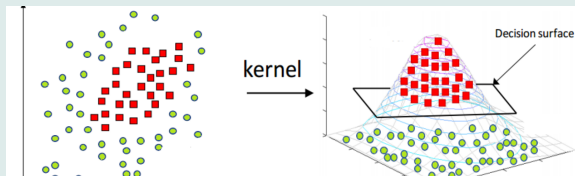
- Correlation: How much variance is explained $\rightarrow \rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y}$
- Covariance: X and Y co-vary. $\rightarrow \text{cov}(X,Y) = \rho_{X,Y} \times \sigma_X \times \sigma_Y$

Altman et al

- Correlation: How much variance is explained $\rightarrow \rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y}$
- Covariance: X and Y co-vary, $\rightarrow \text{cov}(X,Y) = \rho_{X,Y} \sigma_X \sigma_Y$

Altman et al

- Correlation: How much variance is explained $\rightarrow \rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y}$
- Covariance: X and Y co-vary, $\rightarrow \text{cov}(X, Y) = \rho_{X,Y} * \sigma_X * \sigma_Y$



$$k(x_i, x_j) = \varphi(x_i)^T \varphi(x_j) \quad (3)$$

$$\varphi : x \rightarrow \varphi(x) \quad (4)$$

- CCO computes a measure of the conditional dependency for random variables.
- (H_X, K_X) and (H_Y, K_Y) the reproducible kernel Hilbert space (RKHSs) of functions of X and Y respectively.
- (X, Y) a random array on $(X \times Y)$ with distribution $P(X, Y)$
- The cross-covariance operator associated with the pair (X, Y) is the mapping $\Sigma_{X, Y} : H_X \rightarrow H_Y$
- s.t: $E_{X, Y}\{(f(X) - E_X[f(X)])(g(Y) - E_Y[g(Y)])\} : \forall g \in H_Y, f \in H_X$

- CCO computes a measure of the conditional dependency for random variables.
- (H_X, K_X) and (H_Y, K_Y) the reproducible kernel Hilbert space (RKHSs) of functions of X and Y respectively.
- (X, Y) a random array on $(X \times Y)$ with distribution $P(X, Y)$
- The cross-covariance operator associated with the pair (X, Y) is the mapping $\Sigma_{X, Y} : H_X \rightarrow H_Y$
- s.t: $E_{X, Y} \{ (f(X) - E_X[f(X)])(g(Y) - E_Y[g(Y)]) \} :$
 $\forall g \in H_Y, f \in H_X$

- CCO computes a measure of the conditional dependency for random variables.
- (H_X, K_X) and (H_Y, K_Y) the reproducible kernel Hilbert space (RKHSs) of functions of X and Y respectively.
- (X, Y) a random array on $(X \times Y)$ with distribution $P(X, Y)$
- The cross-covariance operator associated with the pair (X, Y) is the mapping $\Sigma_{X, Y} : H_X \rightarrow H_Y$
- s.t: $E_{X, Y} \{ (f(X) - E_X[f(X)])(g(Y) - E_Y[g(Y)]) \} :$
 $\forall g \in H_Y, f \in H_X$

- CCO computes a measure of the conditional dependency for random variables.
- (H_X, K_X) and (H_Y, K_Y) the reproducible kernel Hilbert space (RKHSs) of functions of X and Y respectively.
- (X, Y) a random array on $(X \times Y)$ with distribution $P(X, Y)$
- The cross-covariance operator associated with the pair (X, Y) is the mapping $\Sigma_{X, Y} : H_X \rightarrow H_Y$
- s.t: $E_{X, Y} \{ (f(X) - E_X[f(X)])(g(Y) - E_Y[g(Y)]) \} :$
 $\forall g \in H_Y, f \in H_X$

- CCO computes a measure of the conditional dependency for random variables.
- (H_X, K_X) and (H_Y, K_Y) the reproducible kernel Hilbert space (RKHSs) of functions of X and Y respectively.
- (X, Y) a random array on $(X \times Y)$ with distribution $P(X, Y)$
- The cross-covariance operator associated with the pair (X, Y) is the mapping $\sum_{X, Y} : H_X \rightarrow H_Y$
- s.t: $E_{X, Y} \{ (f(X) - E_X[f(X)])(g(Y) - E_Y[g(Y)]) \} :$
 $\forall g \in H_Y, f \in H_X$

Conditional Covariance Operator (CCO)

- There exists unique bounded operator $V_{Y,X}$, s.t:
- $\langle g, \sum_{YX} f \rangle_{H_Y} = \sum_{YX} = (\sum_{YY})^{1/2} V_{YX} (\sum_{XX})^{1/2}$
- CCO: $\sum_{XX|Y} = \sum_{YY} - (\sum_{YY})^{1/2} V_{YX} V_{XY} (\sum_{XX})^{1/2}$
- CCO captures the conditional variance of Y given X
- $L^2(P_X)$ is the space of all square-integrable¹ functions on X
- If $H_X + R$ is dense in $L^2(P_X)$
- $\langle g, \sum_{XX|X} g \rangle_{H_Y} = E[\text{var}_{Y|X}[g(Y)|X]], \forall g \in H_Y$
- The residual error of $g(Y)$ (where Y is part of H_Y) can be characterized by the CCO
- $\langle g, \sum_{YY|X} g \rangle_{H_Y} = \inf_{f \in H_X} E_{X,Y} \{ (g(Y) - E_Y[g(Y)]) - (f(X) - E_X[f(X)]) \}$

Not intended to be understood at a single slide. Check reference for further understanding.

¹ $\int_{-\infty}^{+\infty} |f(x)|^2 dx < \infty$

Conditional Covariance Operator (CCO)

- There exists unique bounded operator $V_{Y,X}$, s.t:
- $\langle g, \sum_{YX} f \rangle_{H_Y} = \sum_{YX} = (\sum_{YY})^{1/2} V_{YX} (\sum_{XX})^{1/2}$
- CCO: $\sum_{XX|Y} = \sum_{YY} - (\sum_{YY})^{1/2} V_{YX} V_{XY} (\sum_{XX})^{1/2}$
- CCO captures the conditional variance of Y given X
- $L^2(P_X)$ is the space of all square-integrable¹ functions on X
- If $H_X + R$ is dense in $L^2(P_X)$
- $\langle g, \sum_{XX|X} g \rangle_{H_Y} = E[\text{var}_{Y|X}[g(Y)|X]], \forall g \in H_Y$
- The residual error of $g(Y)$ (where Y is part of H_Y) can be characterized by the CCO
- $\langle g, \sum_{YY|X} g \rangle_{H_Y} = \inf_{f \in H_X} E_{X,Y} \{ (g(Y) - E_Y[g(Y)]) - (f(X) - E_X[f(X)]) \}$

Not intended to be understood at a single slide. Check reference for further understanding.

¹ $\int_{-\infty}^{+\infty} |f(x)|^2 dx < \infty$

Conditional Covariance Operator (CCO)

- There exists unique bounded operator $V_{Y,X}$, s.t:
- $\langle g, \sum_{YX} f \rangle_{H_Y} = \sum_{YX} = (\sum_{YY})^{1/2} V_{YX} (\sum_{XX})^{1/2}$
- CCO: $\sum_{XX|Y} = \sum_{YY} - (\sum_{YY})^{1/2} V_{YX} V_{XY} (\sum_{XX})^{1/2}$
- CCO captures the conditional variance of Y given X
- $L^2(P_X)$ is the space of all square-integrable¹ functions on X
- If $H_X + R$ is dense in $L^2(P_X)$
- $\langle g, \sum_{XX|X} g \rangle_{H_Y} = E[\text{var}_{Y|X}[g(Y)|X]], \forall g \in H_Y$
- The residual error of $g(Y)$ (where Y is part of H_Y) can be characterized by the CCO
- $\langle g, \sum_{YY|X} g \rangle_{H_Y} = \inf_{f \in H_X} E_{X,Y} \{ (g(Y) - E_Y[g(Y)]) - (f(X) - E_X[f(X)]) \}$

Not intended to be understood at a single slide. Check reference for further understanding.

¹ $\int_{-\infty}^{+\infty} |f(x)|^2 dx < \infty$

Conditional Covariance Operator(CCO)

- There exists unique bounded operator $V_{Y,X}$, s.t:
- $\langle g, \sum_{YX} f \rangle_{H_Y} = \sum_{YX} = (\sum_{YY})^{1/2} V_{YX} (\sum_{XX})^{1/2}$
- CCO: $\sum_{XX|Y} = \sum_{YY} - (\sum_{YY})^{1/2} V_{YX} V_{XY} (\sum_{XX})^{1/2}$
- CCO captures the conditional variance of Y given X
- $L^2(P_X)$ is the space of all square-integrable¹ functions on X
- If $H_X + R$ is dense in $L^2(P_X)$
- $\langle g, \sum_{XX|X} g \rangle_{H_Y} = E[\text{var}_{Y|X}[g(Y)|X]], \forall g \in H_Y$
- The residual error of $g(Y)$ (where Y is part of H_Y) can be characterized by the CCO
- $\langle g, \sum_{YY|X} g \rangle_{H_Y} = \inf_{f \in H_X} E_{X,Y} \{ (g(Y) - E_Y[g(Y)]) - (f(X) - E_X[f(X)]) \}$

Not intended to be understood at a single slide. Check reference for further understanding.

¹ $\int_{-\infty}^{+\infty} |f(x)|^2 dx < \infty$

Conditional Covariance Operator (CCO)

- There exists unique bounded operator $V_{Y,X}$, s.t:
- $\langle g, \sum_{YX} f \rangle_{H_Y} = \sum_{YX} = (\sum_{YY})^{1/2} V_{YX} (\sum_{XX})^{1/2}$
- CCO: $\sum_{XX|Y} = \sum_{YY} - (\sum_{YY})^{1/2} V_{YX} V_{XY} (\sum_{XX})^{1/2}$
- CCO captures the conditional variance of Y given X
- $L^2(P_X)$ is the space of all square-integrable¹ functions on X
- If $H_X + R$ is dense in $L^2(P_X)$
- $\langle g, \sum_{XX|X} g \rangle_{H_Y} = E[\text{var}_{Y|X}[g(Y)|X]], \forall g \in H_Y$
- The residual error of $g(Y)$ (where Y is part of H_Y) can be characterized by the CCO
- $\langle g, \sum_{YY|X} g \rangle_{H_Y} = \inf_{f \in H_X} E_{X,Y} \{ (g(Y) - E_Y[g(Y)]) - (f(X) - E_X[f(X)]) \}$

Not intended to be understood at a single slide. Check reference for further understanding.

¹ $\int_{-\infty}^{+\infty} |f(x)|^2 dx < \infty$

Conditional Covariance Operator (CCO)

- There exists unique bounded operator $V_{Y,X}$, s.t:
- $\langle g, \sum_{YX} f \rangle_{H_Y} = \sum_{YX} = (\sum_{YY})^{1/2} V_{YX} (\sum_{XX})^{1/2}$
- CCO: $\sum_{XX|Y} = \sum_{YY} - (\sum_{YY})^{1/2} V_{YX} V_{XY} (\sum_{XX})^{1/2}$
- CCO captures the conditional variance of Y given X
- $L^2(P_X)$ is the space of all square-integrable¹ functions on X
- If $H_X + R$ is dense in $L^2(P_X)$
- $\langle g, \sum_{XX|X} g \rangle_{H_Y} = E[\text{var}_{Y|X}[g(Y)|X]], \forall g \in H_Y$
- The residual error of $g(Y)$ (where Y is part of H_Y) can be characterized by the CCO
- $\langle g, \sum_{YY|X} g \rangle_{H_Y} = \inf_{f \in H_X} E_{X,Y} \{ (g(Y) - E_Y[g(Y)]) - (f(X) - E_X[f(X)]) \}$

Not intended to be understood at a single slide. Check reference for further understanding.

¹ $\int_{-\infty}^{+\infty} |f(x)|^2 dx < \infty$

Conditional Covariance Operator (CCO)

- There exists unique bounded operator $V_{Y,X}$, s.t:
- $\langle g, \sum_{YX} f \rangle_{H_Y} = \sum_{YX} = (\sum_{YY})^{1/2} V_{YX} (\sum_{XX})^{1/2}$
- CCO: $\sum_{XX|Y} = \sum_{YY} - (\sum_{YY})^{1/2} V_{YX} V_{XY} (\sum_{XX})^{1/2}$
- CCO captures the conditional variance of Y given X
- $L^2(P_X)$ is the space of all square-integrable¹ functions on X
- If $H_X + R$ is dense in $L^2(P_X)$
- $\langle g, \sum_{XX|X} g \rangle_{H_Y} = E[\text{var}_{Y|X}[g(Y)|X]], \forall g \in H_Y$
- The residual error of $g(Y)$ (where Y is part of H_Y) can be characterized by the CCO
- $\langle g, \sum_{YY|X} g \rangle_{H_Y} = \inf_{f \in H_X} E_{X,Y} \{ (g(Y) - E_Y[g(Y)]) - (f(X) - E_X[f(X)]) \}$

Not intended to be understood at a single slide. Check reference for further understanding.

¹ $\int_{-\infty}^{+\infty} |f(x)|^2 dx < \infty$

Conditional Covariance Operator (CCO)

- There exists unique bounded operator $V_{Y,X}$, s.t:
- $\langle g, \sum_{YX} f \rangle_{H_Y} = \sum_{YX} = (\sum_{YY})^{1/2} V_{YX} (\sum_{XX})^{1/2}$
- CCO: $\sum_{XX|Y} = \sum_{YY} - (\sum_{YY})^{1/2} V_{YX} V_{XY} (\sum_{XX})^{1/2}$
- CCO captures the conditional variance of Y given X
- $L^2(P_X)$ is the space of all square-integrable¹ functions on X
- If $H_X + R$ is dense in $L^2(P_X)$
- $\langle g, \sum_{XX|X} g \rangle_{H_Y} = E[\text{var}_{Y|X}[g(Y)|X]], \forall g \in H_Y$
- The residual error of $g(Y)$ (where Y is part of H_Y) can be characterized by the CCO
- $\langle g, \sum_{YY|X} g \rangle_{H_Y} = \inf_{f \in H_X} E_{X,Y} \{ (g(Y) - E_Y[g(Y)]) - (f(X) - E_X[f(X)]) \}$

Not intended to be understood at a single slide. Check reference for further understanding.

¹ $\int_{-\infty}^{+\infty} |f(x)|^2 dx < \infty$

Conditional Covariance Operator(CCO)

- There exists unique bounded operator $V_{Y,X}$, s.t:
- $\langle g, \sum_{YX} f \rangle_{H_Y} = \sum_{YX} = (\sum_{YY})^{1/2} V_{YX} (\sum_{XX})^{1/2}$
- CCO: $\sum_{XX|Y} = \sum_{YY} - (\sum_{YY})^{1/2} V_{YX} V_{XY} (\sum_{XX})^{1/2}$
- CCO captures the conditional variance of Y given X
- $L^2(P_X)$ is the space of all square-integrable¹ functions on X
- If $H_X + R$ is dense in $L^2(P_X)$
- $\langle g, \sum_{XX|X} g \rangle_{H_Y} = E[\text{var}_{Y|X}[g(Y)|X]], \forall g \in H_Y$
- The residual error of $g(Y)$ (where Y is part of H_Y) can be characterized by the CCO
- $\langle g, \sum_{YY|X} g \rangle_{H_Y} = \inf_{f \in H_X} E_{X,Y} \{ (g(Y) - E_Y[g(Y)]) - (f(X) - E_X[f(X)]) \}$

Not intended to be understood at a single slide. Check reference for further understanding.

¹ $\int_{-\infty}^{+\infty} |f(x)|^2 dx < \infty$

Proposed method

- Let (H_1, k_1) be the RKHS $X \subset R^d$
- Let $T \subseteq [d]$ be a subset of features with cardinality $m \leq d$
- We define $k_1^T(x, \tilde{x}) = k_1(x^T, \tilde{x}^T) \forall x, \tilde{x} \in X$
- k_1 is permutation(π) invariance
 $\forall x, \tilde{x} \in X, k_1(x, \tilde{x}) = k_1(x_\pi, \tilde{x}_\pi)^2$
- $\text{trace}[\sum_{X \times X | Y}]^3$ interpreted as a dependency measure.
- (H, k) is characteristic if $P \rightarrow E_P[k(X, :)]$ is 1to1 map.
- If k is bounded $\rightarrow H + R$ is dense in $L^2(P), \forall P$.

² $(x_{\pi_1}, x_{\pi_2}, \dots, x_{\pi_d})$ as x_π

³ $\text{trace}[A_{(N \times N)}] = \sum_{i=1}^N A_{(i,i)}$

Proposed method

- Let (H_1, k_1) be the RKHS $X \subset R^d$
- Let $T \subseteq [d]$ be a subset of features with cardinality $m \leq d$
- We define $k_1^T(x, \tilde{x}) = k_1(x^T, \tilde{x}^T) \forall x, \tilde{x} \in X$
- k_1 is permutation(π) invariance
 $\forall x, \tilde{x} \in X, k_1(x, \tilde{x}) = k_1(x_\pi, \tilde{x}_\pi)^2$
- $\text{trace}[\sum_{X \times X | Y}]^3$ interpreted as a dependency measure.
- (H, k) is characteristic if $P \rightarrow E_P[k(X, :)]$ is 1to1 map.
- If k is bounded $\rightarrow H + R$ is dense in $L^2(P), \forall P$.

² $(x_{\pi_1}, x_{\pi_2}, \dots, x_{\pi_d})$ as x_π

³ $\text{trace}[A_{(N \times N)}] = \sum_{i=1}^N A_{(i,i)}$

Proposed method

- Let (H_1, k_1) be the RKHS $X \subset R^d$
- Let $T \subseteq [d]$ be a subset of features with cardinality $m \leq d$
- We define $k_1^T(x, \tilde{x}) = k_1(x^T, \tilde{x}^T) \forall x, \tilde{x} \in X$
- k_1 is permutation(π) invariance
 $\forall x, \tilde{x} \in X, k_1(x, \tilde{x}) = k_1(x_\pi, \tilde{x}_\pi)^2$
- $\text{trace}[\sum_{X \times X | Y}]^3$ interpreted as a dependency measure.
- (H, k) is characteristic if $P \rightarrow E_P[k(X, :)]$ is 1to1 map.
- If k is bounded $\rightarrow H + R$ is dense in $L^2(P), \forall P$.

$$^2(x_{\pi_1}, x_{\pi_2}, \dots, x_{\pi_d}) \text{ as } x_\pi$$

$$^3 \text{trace}[A_{(N \times N)}] = \sum_{i=1}^N A_{(i,i)}$$

Proposed method

- Let (H_1, k_1) be the RKHS $X \subset R^d$
- Let $T \subseteq [d]$ be a subset of features with cardinality $m \leq d$
- We define $k_1^T(x, \tilde{x}) = k_1(x^T, \tilde{x}^T) \forall x, \tilde{x} \in X$
- k_1 is permutation(π) invariance
 $\forall x, \tilde{x} \in X, k_1(x, \tilde{x}) = k_1(x_\pi, \tilde{x}_\pi)^2$
 - $\text{trace}[\sum_{X \times Y}]^3$ interpreted as a dependency measure.
 - (H, k) is characteristic if $P \rightarrow E_P[k(X, :)]$ is 1to1 map.
 - If k is bounded $\rightarrow H + R$ is dense in $L^2(P), \forall P$.

$$^2(x_{\pi_1}, x_{\pi_2}, \dots, x_{\pi_d}) \text{ as } x_\pi$$

$$^3 \text{trace}[A_{(N \times m)}] = \sum_{i=1}^N A_{(i,i)}$$

Proposed method

- Let (H_1, k_1) be the RKHS $X \subset R^d$
- Let $T \subseteq [d]$ be a subset of features with cardinality $m \leq d$
- We define $k_1^T(x, \tilde{x}) = k_1(x^T, \tilde{x}^T) \forall x, \tilde{x} \in X$
- k_1 is permutation(π) invariance
 $\forall x, \tilde{x} \in X, k_1(x, \tilde{x}) = k_1(x_\pi, \tilde{x}_\pi)^2$
- $\text{trace}[\sum_{X \times X} Y]^3$ interpreted as a dependency measure.
 - (H, k) is characteristic if $P \rightarrow E_P[k(X, :)]$ is 1to1 map.
 - If k is bounded $\rightarrow H + R$ is dense in $L^2(P), \forall P$.

² $(x_{\pi_1}, x_{\pi_2}, \dots, x_{\pi_d})$ as, x_π
³ $\text{trace}[A_{(N \times N)}] = \sum_{i=1}^N A_{(i,i)}$

Proposed method

- Let (H_1, k_1) be the RKHS $X \subset R^d$
- Let $T \subseteq [d]$ be a subset of features with cardinality $m \leq d$
- We define $k_1^T(x, \tilde{x}) = k_1(x^T, \tilde{x}^T) \forall x, \tilde{x} \in X$
- k_1 is permutation(π) invariance
 $\forall x, \tilde{x} \in X, k_1(x, \tilde{x}) = k_1(x_\pi, \tilde{x}_\pi)^2$
- $\text{trace}[\sum_{X \times X} Y]^3$ interpreted as a dependency measure.
- (H, k) is characteristic if $P \rightarrow E_P[k(X, :)]$ is 1to1 map.
- If k is bounded $\rightarrow H + R$ is dense in $L^2(P), \forall P$.

² $(x_{\pi_1}, x_{\pi_2}, \dots, x_{\pi_d})$ as, x_π

³ $\text{trace}[A_{(N \times N)}] = \sum_{i=1}^N A_{(i,i)}$

Proposed method

- Let (H_1, k_1) be the RKHS $X \subset R^d$
- Let $T \subseteq [d]$ be a subset of features with cardinality $m \leq d$
- We define $k_1^T(x, \tilde{x}) = k_1(x^T, \tilde{x}^T) \forall x, \tilde{x} \in X$
- k_1 is permutation(π) invariance
 $\forall x, \tilde{x} \in X, k_1(x, \tilde{x}) = k_1(x_\pi, \tilde{x}_\pi)^2$
- $\text{trace}[\sum_{X \times X | Y}]^3$ interpreted as a dependency measure.
- (H, k) is characteristic if $P \rightarrow E_P[k(X, :)]$ is 1to1 map.
- If k is bounded $\rightarrow H + R$ is dense in $L^2(P), \forall P$.

² $(x_{\pi_1}, x_{\pi_2}, \dots, x_{\pi_d})$ as, x_π

³ $\text{trace}[A_{(N \times N)}] = \sum_{i=1}^N A_{(i,i)}$

Proposed method

- **L1**⁴: if k_1 is bounded and characteristic $\rightarrow \tilde{k}_1$ is characteristic
- **TH2**⁵: if (H_1, k_1) and (H_2, k_2) are characteristic:

$$\sum_{YY|X} \leq \sum_{YY|X_T}$$

$$\sum_{YY|X} = \sum_{YY|X_T} : \text{iff} : Y \perp X|X_T$$
- **C3**⁶: If (H_1, k_1) is characteristic, $\{y \in [0, 1] : \text{where } \sum_i y_i = 1\} \subset R_k$, and (H_2, k_2) includes the identity function on Y , then we have:

$$\text{Tr}(\sum_{YY|X}) \leq \text{Tr}(\sum_{YY|X_T}), \forall T$$

$$\text{Tr}(\sum_{YY|X}) = \text{Tr}(\sum_{YY|X_T}) : \text{iff} : Y \perp X|X_T$$
- Univariate Objective: $\min_{T:|T|=m} Q(T) = \text{Tr}(\sum_{YY|X_T})$

⁴L \rightarrow Lemma

⁵TH \rightarrow Theorem

⁶C \rightarrow Corollari

- **L1**⁴: if k_1 is bounded and characteristic $\rightarrow \tilde{k}_1$ is characteristic
- **TH2**⁵: if (H_1, k_1) and (H_2, k_2) are characteristic:

$$\sum_{YY|X} \leq \sum_{YY|X_T}$$

$$\sum_{YY|X} = \sum_{YY|X_T} : \text{iff} : Y \perp X|X_T$$
- **C3**⁶: If (H_1, k_1) is characteristic,

$$\{y \in [0, 1] : \text{where } \sum_i y_i = 1\} \subset R_k,$$
 and (H_2, k_2) includes the identity function on Y , then we have:

$$\text{Tr}(\sum_{YY|X}) \leq \text{Tr}(\sum_{YY|X_T}), \forall T$$

$$\text{Tr}(\sum_{YY|X}) = \text{Tr}(\sum_{YY|X_T}) : \text{iff} : Y \perp X|X_T$$
- Univariate Objective: $\min_{T:|T|=m} Q(T) = \text{Tr}(\sum_{YY|X_T})$

⁴L \rightarrow Lemma

⁵TH \rightarrow Theorem

⁶C \rightarrow Corollari

Proposed method

- **L1**⁴: if k_1 is bounded and characteristic $\rightarrow \tilde{k}_1$ is characteristic
- **TH2**⁵: if (H_1, k_1) and (H_2, k_2) are characteristic:

$$\sum_{YY|X} \leq \sum_{YY|X_T}$$

$$\sum_{YY|X} = \sum_{YY|X_T} : \text{iff} : Y \perp X|X_T$$
- **C3**⁶: If (H_1, k_1) is characteristic, $\{y \in [0, 1] : \text{where } \sum_i y_i = 1\} \subset R_k$, and (H_2, k_2) includes the identity function on Y , then we have:

$$\text{Tr}(\sum_{YY|X}) \leq \text{Tr}(\sum_{YY|X_T}), \forall T$$

$$\text{Tr}(\sum_{YY|X}) = \text{Tr}(\sum_{YY|X_T}) : \text{iff} : Y \perp X|X_T$$
- Univariate Objective: $\min_{T:|T|=m} Q(T) = \text{Tr}(\sum_{YY|X_T})$

⁴L \rightarrow Lemma

⁵TH \rightarrow Theorem

⁶C \rightarrow Corollari

Proposed method

- **L1**⁴: if k_1 is bounded and characteristic $\rightarrow \tilde{k}_1$ is characteristic
- **TH2**⁵: if (H_1, k_1) and (H_2, k_2) are characteristic:

$$\sum_{Y|X} \leq \sum_{Y|X_T}$$

$$\sum_{Y|X} = \sum_{Y|X_T} : \text{iff} : Y \perp X|X_T$$
- **C3**⁶: If (H_1, k_1) is characteristic, $\{y \in [0, 1] : \text{where } \sum_i y_i = 1\} \subset R_k$, and (H_2, k_2) includes the identity function on Y , then we have:

$$\text{Tr}(\sum_{Y|X}) \leq \text{Tr}(\sum_{Y|X_T}), \forall T$$

$$\text{Tr}(\sum_{Y|X}) = \text{Tr}(\sum_{Y|X_T}) : \text{iff} : Y \perp X|X_T$$
- Univariate Objective: $\min_{T:|T|=m} Q(T) = \text{Tr}(\sum_{Y|X_T})$

⁴L \rightarrow Lemma

⁵TH \rightarrow Theorem

⁶C \rightarrow Corollari

Proposed method

- **L1**⁴: if k_1 is bounded and characteristic $\rightarrow \tilde{k}_1$ is characteristic
- **TH2**⁵: if (H_1, k_1) and (H_2, k_2) are characteristic:

$$\sum_{Y|X} \leq \sum_{Y|X_T}$$

$$\sum_{Y|X} = \sum_{Y|X_T} : \text{iff} : Y \perp X|X_T$$
- **C3**⁶: If (H_1, k_1) is characteristic, $\{y \in [0, 1] : \text{where } \sum_i y_i = 1\} \subset R_k$, and (H_2, k_2) includes the identity function on Y , then we have:

$$\text{Tr}(\sum_{Y|X}) \leq \text{Tr}(\sum_{Y|X_T}), \forall T$$

$$\text{Tr}(\sum_{Y|X}) = \text{Tr}(\sum_{Y|X_T}) : \text{iff} : Y \perp X|X_T$$
- Univariate Objective: $\min_{T:|T|=m} Q(T) = \text{Tr}(\sum_{Y|X_T})$

⁴L \rightarrow Lemma

⁵TH \rightarrow Theorem

⁶C \rightarrow Corollari

Proposed method

- **L1⁴**: if k_1 is bounded and characteristic $\rightarrow \tilde{k}_1$ is characteristic
- **TH2⁵**: if (H_1, k_1) and (H_2, k_2) are characteristic:

$$\sum_{Y|X} \leq \sum_{Y|X_T}$$

$$\sum_{Y|X} = \sum_{Y|X_T} : \text{iff} : Y \perp X|X_T$$
- **C3⁶**: If (H_1, k_1) is characteristic, $\{y \in [0, 1] : \text{where } \sum_i y_i = 1\} \subset R_k$, and (H_2, k_2) includes the identity function on Y , then we have:

$$\text{Tr}(\sum_{Y|X}) \leq \text{Tr}(\sum_{Y|X_T}), \forall T$$

$$\text{Tr}(\sum_{Y|X}) = \text{Tr}(\sum_{Y|X_T}) : \text{iff} : Y \perp X|X_T$$
- **Univariate Objective**: $\min_{T:|T|=m} Q(T) = \text{Tr}(\sum_{Y|X_T})$

⁴L \rightarrow Lemma

⁵TH \rightarrow Theorem

⁶C \rightarrow Corollari

- **C4:** Let: $\sum_{YY|X_T}$ denote CCO of (X_T, Y) in $(\tilde{H}_1, \tilde{k}_1)$
 denote: $F_m = \tilde{H}_1 + R = f + c : f \in \tilde{H}_1, c \in R$
 then: $Tr(\sum_{YY|X_T}) = \epsilon_{F_m}(X_T) = \inf_{f \in F} E_{X,Y}(Y - f(X_T))^2$

- Given n samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the empirical estimate is given by:

$$Tr(\sum_{YY|X_T}^{(n)}) = \text{trace} \left\{ \sum_{YY}^{(n)} - \sum_{Y, X_T}^{(n)} [\sum_{X_T, Y_T}^{(n)} + \epsilon I] \sum_{X_T Y}^{(n)} \right\}$$

$$Tr(\sum_{YY|X_T}^{(n)}) = \epsilon_n \text{trace} \{ G_Y [G_X + n\epsilon_n I_n] \}$$

$$\text{where } G_X = (I_n - \frac{1}{n} I * I^T) K_{X_T} (I_n - \frac{1}{n} I * I^T)$$

$$\text{and } G_Y = (I_n - \frac{1}{n} I * I^T) K_{Y_T} (I_n - \frac{1}{n} I * I^T)$$

- **WLG:** $\text{trace}[G_Y(G_{X_T} + n\epsilon_n I_n)^{-1}] = \text{trace}[Y Y^T (G_{X_T} + n\epsilon_n I_n)^{-1}] = \text{trace}[Y^T (G_{X_T} + n\epsilon_n I_n)^{-1} Y]$
- **Univariate Objective:** $\min_{|T|=m} Q^{(n)} = Y^T (G_{X_T} + n\epsilon_n I_n)^{-1} Y$
 where $Y = (y_1, y_2, \dots, y_n)$ is a n -dimensional vector.

Not intended to be understood at a single slide. Check reference for further understanding.

- **C4:** Let: $\sum_{YY|X_T}$ denote CCO of (X_T, Y) in $(\tilde{H}_1, \tilde{k}_1)$
 denote: $F_m = \tilde{H}_1 + R = f + c : f \in \tilde{H}_1, c \in R$
 then: $Tr(\sum_{YY|X_T}) = \epsilon_{F_m}(X_T) = \inf_{f \in F} E_{X,Y}(Y - f(X_T))^2$

- Given n samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the empirical estimate is given by:

$$Tr(\sum_{YY|X_T}^{(n)}) = \text{trace} \left\{ \sum_{YY}^{(n)} - \sum_{Y, X_T}^{(n)} [\sum_{X_T, Y_T}^{(n)} + \epsilon I] \sum_{X_T Y}^{(n)} \right\}$$

$$Tr(\sum_{YY|X_T}^{(n)}) = \epsilon_n \text{trace} \{ G_Y [G_X + n\epsilon_n I_n] \}$$

$$\text{where } G_X = (I_n - \frac{1}{n} I * I^T) K_{X_T} (I_n - \frac{1}{n} I * I^T)$$

$$\text{and } G_Y = (I_n - \frac{1}{n} I * I^T) K_{Y_T} (I_n - \frac{1}{n} I * I^T)$$

- **WLG:** $\text{trace}[G_Y(G_{X_T} + n\epsilon_n I_n)^{-1}] = \text{trace}[Y Y^T (G_{X_T} + n\epsilon_n I_n)^{-1}] = \text{trace}[Y^T (G_{X_T} + n\epsilon_n I_n)^{-1} Y]$
- **Univariate Objective:** $\min_{|T|=m} Q^{(n)} = Y^T (G_{X_T} + n\epsilon_n I_n)^{-1} Y$
 where $Y = (y_1, y_2, \dots, y_n)$ is a n -dimensional vector.

Not intended to be understood at a single slide. Check reference for further understanding.

- **C4:** Let: $\sum_{YY|X_T}$ denote CCO of (X_T, Y) in $(\tilde{H}_1, \tilde{k}_1)$
 denote: $F_m = \tilde{H}_1 + R = f + c : f \in \tilde{H}_1, c \in R$
 then: $Tr(\sum_{YY|X_T}) = \epsilon_{F_m}(X_T) = \inf_{f \in F} E_{X,Y}(Y - f(X_T))^2$

- Given n samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the empirical estimate is given by:

$$Tr(\sum_{YY|X_T}^{(n)}) = \text{trace} \left\{ \sum_{YY}^{(n)} - \sum_{Y, X_T}^{(n)} [\sum_{X_T, Y_T}^{(n)} + \epsilon I] \sum_{X_T Y}^{(n)} \right\}$$

$$Tr(\sum_{YY|X_T}^{(n)}) = \epsilon_n \text{trace} \{ G_Y [G_X + n\epsilon_n I_n] \}$$

$$\text{where } G_X = (I_n - \frac{1}{n} I * I^T) K_{X_T} (I_n - \frac{1}{n} I * I^T)$$

$$\text{and } G_Y = (I_n - \frac{1}{n} I * I^T) K_{Y_T} (I_n - \frac{1}{n} I * I^T)$$

- WLG: $\text{trace}[G_Y(G_{X_T} + n\epsilon_n I_n)^{-1}] = \text{trace}[Y Y^T (G_{X_T} + n\epsilon_n I_n)^{-1}] = \text{trace}[Y^T (G_{X_T} + n\epsilon_n I_n)^{-1} Y]$
- Univariate Objective: $\min_{|T|=m} Q^{(n)} = Y^T (G_{X_T} + n\epsilon_n I_n)^{-1} Y$
 where $Y = (y_1, y_2, \dots, y_n)$ is a n -dimensional vector.

Not intended to be understood at a single slide. Check reference for further understanding.

Proposed method

- **C4:** Let: $\sum_{YY|X_T}$ denote CCO of (X_T, Y) in $(\tilde{H}_1, \tilde{k}_1)$
 denote: $F_m = \tilde{H}_1 + R = f + c : f \in \tilde{H}_1, c \in R$
 then: $Tr(\sum_{YY|X_T}) = \epsilon_{F_m}(X_T) = \inf_{f \in F} E_{X,Y}(Y - f(X_T))^2$
- Given n samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the empirical estimate is given by:

$$Tr(\sum_{YY|X_T}^{(n)}) = \text{trace} \left\{ \sum_{YY}^{(n)} - \sum_{Y, X_T}^{(n)} [\sum_{X_T, Y_T}^{(n)} + \epsilon I] \sum_{X_T Y}^{(n)} \right\}$$

$$Tr(\sum_{YY|X_T}^{(n)}) = \epsilon_n \text{trace} \{ G_Y [G_X + n \epsilon_n I_n] \}$$

$$\text{where } G_X = (I_n - \frac{1}{n} I * I^T) K_{X_T} (I_n - \frac{1}{n} I * I^T)$$

$$\text{and } G_Y = (I_n - \frac{1}{n} I * I^T) K_{Y_T} (I_n - \frac{1}{n} I * I^T)$$

- WLG: $\text{trace}[G_Y (G_{X_T} + n \epsilon_n I_n)^{-1}] = \text{trace}[Y Y^T (G_{X_T} + n \epsilon_n I_n)^{-1}] = \text{trace}[Y^T (G_{X_T} + n \epsilon_n I_n)^{-1} Y]$
- Univariate Objective: $\min_{|T|=m} Q^{(n)} = Y^T (G_{X_T} + n \epsilon_n I_n)^{-1} Y$
 where $Y = (y_1, y_2, \dots, y_n)$ is a n -dimensional vector.

Not intended to be understood at a single slide. Check reference for further understanding.

Proposed method

- **C4:** Let: $\sum_{YY|X_T}$ denote CCO of (X_T, Y) in $(\tilde{H}_1, \tilde{k}_1)$
 denote: $F_m = \tilde{H}_1 + R = f + c : f \in \tilde{H}_1, c \in R$
 then: $Tr(\sum_{YY|X_T}) = \epsilon_{F_m}(X_T) = \inf_{f \in F} E_{X,Y}(Y - f(X_T))^2$
- Given n samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the empirical estimate is given by:

$$Tr(\sum_{YY|X_T}^{(n)}) = \text{trace} \left\{ \sum_{YY}^{(n)} - \sum_{Y, X_T}^{(n)} [\sum_{X_T, Y_T}^{(n)} + \epsilon I] \sum_{X_T Y}^{(n)} \right\}$$

$$Tr(\sum_{YY|X_T}^{(n)}) = \epsilon_n \text{trace} \{ G_Y [G_X + n \epsilon_n I_N] \}$$
 where $G_X = (I_n - \frac{1}{n} I * I^T) K_{X_T} (I_n - \frac{1}{n} I * I^T)$
 and $G_Y = (I_n - \frac{1}{n} I * I^T) K_{Y_T} (I_n - \frac{1}{n} I * I^T)$
- WLG: $\text{trace}[G_Y (G_{X_T} + n \epsilon_n I_N)^{-1}] = \text{trace}[Y Y^T (G_{X_T} + n \epsilon_n I_N)^{-1}] = \text{trace}[Y^T (G_{X_T} + n \epsilon_n I_N)^{-1} Y]$
- Univariate Objective: $\min_{|T|=m} Q^{(n)} = Y^T (G_{X_T} + n \epsilon_n I_N)^{-1} Y$
 where $Y = (y_1, y_2, \dots, y_n)$ is a n -dimensional vector.

Not intended to be understood at a single slide. Check reference for further understanding.

- **C4:** Let: $\sum_{YY|X_T}$ denote CCO of (X_T, Y) in $(\tilde{H}_1, \tilde{k}_1)$
 denote: $F_m = \tilde{H}_1 + R = f + c : f \in \tilde{H}_1, c \in R$
 then: $Tr(\sum_{YY|X_T}) = \epsilon_{F_m}(X_T) = \inf_{f \in F} E_{X,Y}(Y - f(X_T))^2$
- Given n samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the empirical estimate is given by:

$$Tr(\sum_{YY|X_T}^{(n)}) = \text{trace} \left\{ \sum_{YY}^{(n)} - \sum_{Y, X_T}^{(n)} [\sum_{X_T, Y_T}^{(n)} + \epsilon I] \sum_{X_T Y}^{(n)} \right\}$$

$$Tr(\sum_{YY|X_T}^{(n)}) = \epsilon_n \text{trace} \{ G_Y [G_X + n \epsilon_n I_n] \}$$
 where $G_X = (I_n - \frac{1}{n} I * I^T) K_{X_T} (I_n - \frac{1}{n} I * I^T)$
 and $G_Y = (I_n - \frac{1}{n} I * I^T) K_{Y_T} (I_n - \frac{1}{n} I * I^T)$
- WLG: $\text{trace}[G_Y (G_{X_T} + n \epsilon_n I_n)^{-1}] = \text{trace}[Y Y^T (G_{X_T} + n \epsilon_n I_n)^{-1}] = \text{trace}[Y^T (G_{X_T} + n \epsilon_n I_n)^{-1} Y]$
- Univariate Objective: $\min_{|T|=m} Q^{(n)} = Y^T (G_{X_T} + n \epsilon_n I_n)^{-1} Y$
 where $Y = (y_1, y_2, \dots, y_n)$ is a n -dimensional vector.

Not intended to be understood at a single slide. Check reference for further understanding.

Proposed method

- **C4:** Let: $\sum_{YY|X_T}$ denote CCO of (X_T, Y) in $(\tilde{H}_1, \tilde{k}_1)$
 denote: $F_m = \tilde{H}_1 + R = f + c : f \in \tilde{H}_1, c \in R$
 then: $Tr(\sum_{YY|X_T}) = \epsilon_{F_m}(X_T) = \inf_{f \in F} E_{X,Y}(Y - f(X_T))^2$

- Given n samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the empirical estimate is given by:

$$Tr(\sum_{YY|X_T}^{(n)}) = \text{trace} \left\{ \sum_{YY}^{(n)} - \sum_{Y, X_T}^{(n)} [\sum_{X_T, Y_T}^{(n)} + \epsilon I] \sum_{X_T Y}^{(n)} \right\}$$

$$Tr(\sum_{YY|X_T}^{(n)}) = \epsilon_n \text{trace} \{ G_Y [G_X + n \epsilon_n I_n] \}$$

$$\text{where } G_X = (I_n - \frac{1}{n} I * I^T) K_{X_T} (I_n - \frac{1}{n} I * I^T)$$

$$\text{and } G_Y = (I_n - \frac{1}{n} I * I^T) K_{Y_T} (I_n - \frac{1}{n} I * I^T)$$

- WLG: $\text{trace}[G_Y (G_{X_T} + n \epsilon_n I_n)^{-1}] = \text{trace}[Y Y^T (G_{X_T} + n \epsilon_n I_n)^{-1}] = \text{trace}[Y^T (G_{X_T} + n \epsilon_n I_n)^{-1} Y]$
- Univariate Objective: $\min_{|T|=m} Q^{(n)} = Y^T (G_{X_T} + n \epsilon_n I_n)^{-1} Y$
 where $Y = (y_1, y_2, \dots, y_n)$ is a n -dimensional vector.

Not intended to be understood at a single slide. Check reference for further understanding.

Proposed method

- **C4:** Let: $\sum_{YY|X_T}$ denote CCO of (X_T, Y) in $(\tilde{H}_1, \tilde{k}_1)$
denote: $F_m = \tilde{H}_1 + R = f + c : f \in \tilde{H}_1, c \in R$
then: $Tr(\sum_{YY|X_T}) = \epsilon_{F_m}(X_T) = \inf_{f \in F} E_{X,Y}(Y - f(X_T))^2$

- Given n samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the empirical estimate is given by:

$$Tr(\sum_{YY|X_T}^{(n)}) = trace\left\{\sum_{YY}^{(n)} - \sum_{Y, X_T}^{(n)}[\sum_{X_T, Y_T}^{(n)} + \epsilon I]\sum_{X_T Y}^{(n)}\right\}$$

$$Tr(\sum_{YY|X_T}^{(n)}) = \epsilon_n trace\{G_Y[G_X + n\epsilon_n I_n]\}$$

$$\text{where } G_X = (I_n - \frac{1}{n}I * I^T)K_{X_T}(I_n - \frac{1}{n}I * I^T)$$

$$\text{and } G_Y = (I_n - \frac{1}{n}I * I^T)K_{Y_T}(I_n - \frac{1}{n}I * I^T)$$

- **WLG:** $trace[G_Y(G_{X_T} + n\epsilon_n I_n)^{-1}] = trace[Y Y^T (G_{X_T} + n\epsilon_n I_n)^{-1}] = trace[Y^T (G_{X_T} + n\epsilon_n I_n)^{-1} Y]$
- **Univariate Objective:** $min_{|T|=m} Q^{(n)} = Y^T (G_{X_T} + n\epsilon_n I_n)^{-1} Y$
where $Y = (y_1, y_2, \dots, y_n)$ is a n -dimensional vector.

Not intended to be understood at a single slide. Check reference for further understanding.

Proposed method

- C4:** Let: $\sum_{YY|X_T}$ denote CCO of (X_T, Y) in $(\tilde{H}_1, \tilde{k}_1)$
 denote: $F_m = \tilde{H}_1 + R = f + c : f \in \tilde{H}_1, c \in R$
 then: $Tr(\sum_{YY|X_T}) = \epsilon_{F_m}(X_T) = \inf_{f \in F} E_{X,Y}(Y - f(X_T))^2$
- Given n samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the empirical estimate is given by:

$$Tr(\sum_{YY|X_T}^{(n)}) = \text{trace} \left\{ \sum_{YY}^{(n)} - \sum_{Y, X_T}^{(n)} (\sum_{X_T, Y_T}^{(n)} + \epsilon I) \sum_{X_T Y}^{(n)} \right\}$$

$$Tr(\sum_{YY|X_T}^{(n)}) = \epsilon_n \text{trace} \{ G_Y [G_X + n \epsilon_n I_n] \}$$

where $G_X = (I_n - \frac{1}{n} I * I^T) K_{X_T} (I_n - \frac{1}{n} I * I^T)$
 and $G_Y = (I_n - \frac{1}{n} I * I^T) K_{Y_T} (I_n - \frac{1}{n} I * I^T)$
- WLG:** $\text{trace}[G_Y(G_{X_T} + n \epsilon_n I_n)^{-1}] = \text{trace}[Y Y^T (G_{X_T} + n \epsilon_n I_n)^{-1}] = \text{trace}[Y^T (G_{X_T} + n \epsilon_n I_n)^{-1} Y]$
- Univariate Objective: $\min_{|T|=m} Q^{(n)} = Y^T (G_{X_T} + n \epsilon_n I_n)^{-1} Y$
 where $Y = (y_1, y_2, \dots, y_n)$ is a n -dimensional vector.

Not intended to be understood at a single slide. Check reference for further understanding.

Proposed method

- **C4:** Let: $\sum_{YY|X_T}$ denote CCO of (X_T, Y) in $(\tilde{H}_1, \tilde{k}_1)$
 denote: $F_m = \tilde{H}_1 + R = f + c : f \in \tilde{H}_1, c \in R$
 then: $Tr(\sum_{YY|X_T}) = \epsilon_{F_m}(X_T) = \inf_{f \in F} E_{X,Y}(Y - f(X_T))^2$
- Given n samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the empirical estimate is given by:

$$Tr(\sum_{YY|X_T}^{(n)}) = \text{trace} \left\{ \sum_{YY}^{(n)} - \sum_{Y, X_T}^{(n)} (\sum_{X_T, Y_T}^{(n)} + \epsilon I) \sum_{X_T Y}^{(n)} \right\}$$

$$Tr(\sum_{YY|X_T}^{(n)}) = \epsilon_n \text{trace} \{ G_Y [G_X + n \epsilon_n I_n] \}$$

$$\text{where } G_X = (I_n - \frac{1}{n} I + I^T) K_{X_T} (I_n - \frac{1}{n} I + I^T)$$

$$\text{and } G_Y = (I_n - \frac{1}{n} I + I^T) K_{Y_T} (I_n - \frac{1}{n} I + I^T)$$

- **WLG:** $\text{trace}[G_Y(G_{X_T} + n \epsilon_n I_n)^{-1}] = \text{trace}[Y Y^T (G_{X_T} + n \epsilon_N I_N)^{-1}] = \text{trace}[Y^T (G_{X_T} + n \epsilon_N I_N)^{-1} Y]$
- **Univariate Objective:** $\min_{|T|=m} Q^{(n)} = Y^T (G_{X_T} + n \epsilon_N I_N)^{-1} Y$
 where $Y = (y_1, y_2, \dots, y_n)$ is a n -dimensional vector.

Not intended to be understood at a single slide. Check reference for further understanding.

- $\operatorname{argmin}_w : Y^T (G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} Y$
 subject to: $w_i \in \{0, 1\}, i = 1, \dots, d$
 where $I^T w \leq m$

$\odot \Rightarrow$ Hadamard

- $\operatorname{argmin}_w : Y^T (G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} Y$
 subject to: $w_i \in \{0, 1\}, i = 1, \dots, d$
 where $I^T w \leq m$

$\odot \Rightarrow$ Hadamard

- $\operatorname{argmin}_w : Y^T (G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} Y$
subject to: $w_i \in \{0, 1\}, i = 1, \dots, d$
where $I^T w \leq m$

$\odot \Rightarrow$ Hadamard

- $\operatorname{argmin}_w : Y^T (G_{X_w \odot X} + n \epsilon_N I_N)^{-1} Y$
subject to : $w_i \in \{0, 1\}, i = 1, \dots, d$
where $1^T w \leq m$

- $\operatorname{argmin}_w : Y^T (G_{X_w \odot X} + n \epsilon_N I_N)^{-1} Y$
subject to : $w_i \in \{0, 1\}, i = 1, \dots, d$
where $I^T w \leq m$

- $\operatorname{argmin}_w : Y^T (G_{X_w \odot X} + n \epsilon_N I_N)^{-1} Y$
subject to : $w_i \in \{0, 1\}, i = 1, \dots, d$
where $I^T w \leq m$

- $\operatorname{argmin}_w : Y^T (G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} Y$
subject to $0 \leq w_i \leq 1, i = 1, \dots, d$
where $I^T w \leq m$

- $\operatorname{argmin}_w : Y^T (G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} Y + \lambda_1 [l^T w - m]$
subject to $0 \leq w_i \leq 1, i = 1, \dots, d$
where $\lambda_1 \geq 0$

- $\text{argmin}_{w, \alpha} : \alpha * y + \|(G_{X_{w \odot X}} + n\epsilon_N I_N)\alpha + y\|_2^2$
 subject to $0 \leq w_i \leq 1, i = 1, \dots, d$
 where $I^T w \leq m$ and $\alpha = (G_{X_{w \odot X}} + n\epsilon_N I_N)^{-1} y$

- $\operatorname{argmin}_w : Y^T (G_{X_w \otimes X} + n \epsilon_N I_N)^{-1} Y$
 subject to $0 \leq w_i \leq 1, i = 1, \dots, d$
 where $I^T w \leq m$

$$(G_{X_w \otimes X} + n \epsilon_N I_N)^{-1} \approx \frac{1}{c_n n} I - \frac{1}{c_n^2 n^2} V (I_D + \frac{1}{c_n n} V_w^T V_w)^{-1} V_w$$

$$(G_{X_w \otimes X} + n \epsilon_N I_N)^{-1} \approx \frac{1}{c_n n} (I - V_w (V_w^T V_w + c_n b I_D)^{-1} V_w^T)$$

- $\operatorname{argmin}_w : Y^T (G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} Y$
 subject to $0 \leq w_i \leq 1, i = 1, \dots, d$
 where $I^T w \leq m$

$$(G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} \approx \frac{1}{\epsilon_n n} I - \frac{1}{\epsilon_n^2 n^2} V (I_D + \frac{1}{\epsilon_n n} V_w^T V_w)^{-1} V_w$$

$$(G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} \approx \frac{1}{\epsilon_n n} (I - V_w (V_w^T V_w + \epsilon_n b I_D)^{-1} V_w^T)$$

- $\operatorname{argmin}_w : Y^T (G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} Y$
 subject to $0 \leq w_i \leq 1, i = 1, \dots, d$
 where $I^T w \leq m$

$$(G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} \approx \frac{1}{\epsilon_n n} I - \frac{1}{\epsilon_n^2 n^2} V (I_D + \frac{1}{\epsilon_n n} V_w^T V_w)^{-1} V_w$$

$$(G_{X_{w \odot X}} + n \epsilon_N I_N)^{-1} \approx \frac{1}{\epsilon_n n} (I - V_w (V_w^T V_w + \epsilon_n b I_D)^{-1} V_w^T)$$

racefox

Longer. Faster. Forever.

