

# Information Theory and Feature Selection

(Joint Informativeness and Tractability)

Leonidas Lefakis

Zalando Research Labs

- ▶ Construction

$$X_1, \dots, X_D \rightarrow f_1(X_1, \dots, X_D), \dots, f_k(X_1, \dots, X_D)$$

- ▶ Construction

$$X_1, \dots, X_D \rightarrow f_1(X_1, \dots, X_D), \dots, f_k(X_1, \dots, X_D)$$

### Examples

- ▶ Principal Component Analysis (PCA)
- ▶ Linear Discriminant Analysis (LDA)
- ▶ Autoencoders (Neural Networks)

- ▶ Selection

$$X_1, \dots, X_D \rightarrow X_{s_1}, \dots, X_{s_k}$$

### Approaches

- ▶ Wrappers
- ▶ Embedded methods
- ▶ Filters

- ▶ Selection

$$X_1, \dots, X_D \rightarrow X_{s_1}, \dots, X_{s_k}$$

- ▶ Wrappers

Features are selected relative to the performance of a specific predictor.

Example [RFE-SVM](#).

- ▶ Selection

$$X_1, \dots, X_D \rightarrow X_{s_1}, \dots, X_{s_k}$$

- ▶ Embedded Methods

Features are selected internally while optimizing the predictor.

Example [Decision Trees](#).

- ▶ Selection

$$X_1, \dots, X_D \rightarrow X_{s_1}, \dots, X_{s_k}$$

- ▶ Filters

Features are assessed based on some goodness-of-fit function  $\Phi$  that is classifier agnostic.

Example [Correlation](#).

- ▶ Entropy

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x)$$

- ▶ Joint Entropy

$$H(X, Y) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y)$$

- ▶ Conditional Entropy

$$H(Y|X) = - \sum_{x \in \mathcal{X}} p(x) H(Y|X = x)$$



- ▶ Relative Entropy (Kullback-Leibler Divergence )

$$D(p||q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)}$$

- ▶ Mutual Information

$$I(X; Y) = \underbrace{\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right) dx dy}_{D_{KL}(p(x, y) || p(x)p(y))}$$

$$I(X; Y) = \underbrace{\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right) dx dy}_{D_{KL}(p(x, y) || p(x)p(y))}$$

$$I(X; Y) = H(Y) - H(Y|X)$$

Reduction in uncertainty of  $Y$  if  $X$  is known

$$X \perp\!\!\!\perp Y \Rightarrow I(X; Y) = 0$$

$$Y = f(X) \Rightarrow I(X; Y) = H(Y)$$

We will look into feature selection in the context of classification.  
Given features

$$\{X_1, \dots, X_D\} \in \mathbb{R}$$

and a class variable  $Y$ , we wish to select a subset  $S$  of size  $K \ll D$  such that a predictor

$$f : \mathbb{R}^K \rightarrow Y$$

trained in this projected subspace generalizes well.

Given

$$X \times Y \in \mathbb{R}^D \times \{1 \dots C\}, \quad f(X) = \hat{Y}$$

We define the error variable

$$E = \begin{cases} 0 & \text{if } \hat{Y} = Y \\ 1 & \text{if } \hat{Y} \neq Y \end{cases}$$

$$\begin{aligned} H(E, Y | \hat{Y}) &\stackrel{(1)}{=} H(Y | \hat{Y}) + \underbrace{H(E | Y, \hat{Y})}_{=0} \\ &\stackrel{(1)}{=} \underbrace{H(E | \hat{Y})}_{\leq H(E)} + H(Y | E, \hat{Y}) \end{aligned}$$

$$H(A, B) = H(A) + H(B|A) \tag{1}$$

$$\begin{aligned}
 H(Y|\hat{Y}) &= H(E, Y|\hat{Y}) \\
 &\leq H(E) + H(Y|E, \hat{Y}) \\
 &\leq 1 + P_e \log(|\mathcal{Y}|-1)
 \end{aligned}$$

$$H(E) = H(B(1, P_e)) \leq H\left(B\left(1, \frac{1}{2}\right)\right) = 1$$

$$H(Y|E, \hat{Y}) = (1 - P_e) \underbrace{H(Y|E=0, \hat{Y})}_{=0} + P_e \underbrace{H(Y|E=1, \hat{Y})}_{\leq \log(|\mathcal{Y}|-1)}$$

$$H(Y|\hat{Y}) \leq 1 + P_e \log(|\mathcal{Y}|-1)$$

$$\stackrel{(2)}{\implies} H(Y) - I(Y; \hat{Y}) \leq 1 + P_e \log(|\mathcal{Y}|-1)$$

$$\implies P_e \geq \frac{H(Y) - I(Y; \hat{Y}) - 1}{\log(|\mathcal{Y}|-1)}$$

$$\stackrel{3}{\implies} P_e \geq \frac{H(Y) - 1 - I(X; Y)}{\log(|\mathcal{Y}|-1)}$$

$$I(A; B) = H(A) - H(A|B) \tag{2}$$

$$Y \rightarrow X \rightarrow Z \implies I(Y; X) \geq I(Y; Z) \tag{3}$$

$$H(Y|\hat{Y}) \leq 1 + P_e \log(|\mathcal{Y}|-1)$$

$$\stackrel{(2)}{\implies} H(Y) - I(Y; \hat{Y}) \leq 1 + P_e \log(|\mathcal{Y}|-1)$$

$$\implies P_e \geq \frac{H(Y) - I(Y; \hat{Y}) - 1}{\log(|\mathcal{Y}|-1)}$$

$$\stackrel{3}{\implies} P_e \geq \frac{H(Y) - 1 - I(X; Y)}{\log(|\mathcal{Y}|-1)}$$

$$I(A; B) = H(A) - H(A|B) \tag{2}$$

$$Y \rightarrow X \rightarrow Z \implies I(Y; X) \geq I(Y; Z) \tag{3}$$



$$H(Y|\hat{Y}) \leq 1 + P_e \log(|\mathcal{Y}|-1)$$

$$\stackrel{(2)}{\implies} H(Y) - I(Y; \hat{Y}) \leq 1 + P_e \log(|\mathcal{Y}|-1)$$

$$\implies P_e \geq \frac{H(Y) - I(Y; \hat{Y}) - 1}{\log(|\mathcal{Y}|-1)}$$

$$\stackrel{3}{\implies} P_e \geq \frac{H(Y) - 1 - I(X; Y)}{\log(|\mathcal{Y}|-1)}$$

$$I(A; B) = H(A) - H(A|B) \tag{2}$$

$$Y \rightarrow X \rightarrow Z \implies I(Y; X) \geq I(Y; Z) \tag{3}$$

$$H(Y|\hat{Y}) \leq 1 + P_e \log(|\mathcal{Y}|-1)$$

$$\stackrel{(2)}{\implies} H(Y) - I(Y; \hat{Y}) \leq 1 + P_e \log(|\mathcal{Y}|-1)$$

$$\implies P_e \geq \frac{H(Y) - I(Y; \hat{Y}) - 1}{\log(|\mathcal{Y}|-1)}$$

$$\stackrel{3}{\implies} P_e \geq \frac{H(Y) - 1 - I(X; Y)}{\log(|\mathcal{Y}|-1)}$$

$$I(A; B) = H(A) - H(A|B) \tag{2}$$

$$Y \rightarrow X \rightarrow Z \implies I(Y; X) \geq I(Y; Z) \tag{3}$$

$$S = \operatorname{argmax}_{S', |S'|=K} I(X_{S'(1)}, X_{S'(2)}, \dots, X_{S'(K)}; Y)$$

$$S = \operatorname{argmax}_{S', |S'|=K} I(X_{S'(1)}, X_{S'(2)}, \dots, X_{S'(K)}; Y)$$

NP-HARD

$$S = \operatorname{argmax}_{S', |S'|=K} \sum_{k=1}^K I(X_{S'(k)}; Y)$$

- ▶ Considers the Relevance of each variable individually
- ▶ Does not consider Redundancy
- ▶ Does not consider Joint Informativeness

- ▶ **mRMR** : Feature Selection Based on Mutual Information: Criteria of Max-Dependency, Max-Relevance, and Min-Redundancy, H.Peng et al.
- ▶ **CMIM** : Fast Binary Feature Selection with Conditional Mutual Information, F.Fleuret

Relevance

$$D(S, Y) = \frac{1}{|S|} \sum_{X_d \in S} I(X_d; Y)$$

Redundancy

$$R(S) = \frac{1}{|S|^2} \sum_{X_d \in S} \sum_{X_l \in S} I(X_d; X_l)$$

mRMR

$$\Phi(S, Y) = D(S, Y) - R(S)$$

## Forward Selection



```
 $S_0 = \emptyset$   
for  $k = 1 \dots K$  do  
   $z^* = 0$   
  for  $X_j \in F \setminus S_{k-1}$  do  
     $S' \leftarrow S_{k-1} \cup X_j$   
     $z \leftarrow \Phi(S', Y)$   
    if  $z > z^*$  then  
       $s^* \leftarrow s$   
       $S^* \leftarrow S'$   
    end if  
  end for  
   $S_i \leftarrow S^*$   
end for  
return  $S_K$ 
```



- ▶ Discretize
- ▶ Distributions approximated using Parzen windows

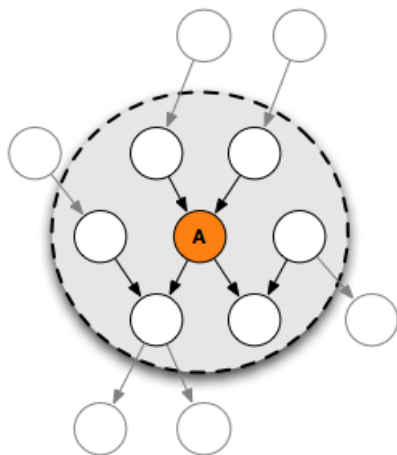
$$\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \delta(x - x^{(i)}, h)$$

$$\delta(dx, h) = \frac{1}{\sqrt{2\pi^D}} \exp\left(-\frac{dx^T \Sigma^{-1} dx}{2h^2}\right)$$

- ▶ Parametric Model

# CMIM (Binary Features)

Markov Blanket



Markov Blanket of variable  $A$ <sup>1</sup>

---

<sup>1</sup>Wiki Commons

# CMIM (Binary Features)

## Markov Blanket



For a set  $M$  of variables that form a Markov Blanket of  $X_i$  we have

$$p(F \setminus \{X_i, M\}, Y | X_i, M) = p(F \setminus \{X_i, M\}, Y | M)$$

$$I(F \setminus \{X_i, M\}, Y; X_i | M) = 0$$

For a set  $M$  of variables that form a Markov Blanket of  $X_i$  we have

$$p(F \setminus \{X_i, M\}, Y | X_i, M) = p(F \setminus \{X_i, M\}, Y | M)$$

$$I(F \setminus \{X_i, M\}, Y; X_i | M) = 0$$

For Feature Selection  $\rightarrow$  Too Strong

$$I(Y; X_i | M) = 0$$

$$I(X_1, \dots, X_K; Y) = H(Y) - \underbrace{H(Y|X_1, \dots, X_K)}_{\text{intractable}}$$

$$I(X_1, \dots, X_K; Y) = H(Y) - \underbrace{H(Y|X_1, \dots, X_K)}_{\text{intractable}}$$

$$I(X_i; Y|X_j) = H(Y|X_j) - H(Y|X_i, X_j)$$

Distributions over triplets of variables

$$I(X_1, \dots, X_K; Y) = H(Y) - \underbrace{H(Y|X_1, \dots, X_K)}_{\text{intractable}}$$

$$I(X_i; Y|X_j) = H(Y|X_j) - H(Y|X_i, X_j)$$

$$S_1 \leftarrow \operatorname{argmax}_d \hat{I}(X_d; Y)$$

$$S_k \leftarrow \operatorname{argmax}_d \left\{ \min_{l < k} \hat{I}(Y; X_d | X_{S_l}) \right\}$$

Distributions over triplets of variables

$$S_k \leftarrow \operatorname{argmax}_d \left\{ \underbrace{\min_{l < k} \hat{I}(Y; X_d | X_{S_l})}_{\text{Can only decrease}} \right\}$$



$$S_k \leftarrow \operatorname{argmax}_d \{ \underbrace{\min_{l < k} \hat{I}(Y; X_d | X_{S_l})}_{\text{Can only decrease}} \}$$

	1	2	3	4	5	6	7
1	6	3	1	5	4	2	5
2	3	5	?	5	3	?	3
3	5	2	?	4	?	?	?
4	4	?	?	6	?	?	?
5	3	?	?	4	?	?	?
ps[n]	3	2	1	4	3	2	3
m[n]	5	3	1	5	2	1	2

$$ps[n] = \min_{l < m(n)} \hat{I}(Y; X_d | X_{S_l})$$

$$S = \operatorname{argmax}_{S', |S'|=K} I(X_{S'(1)}, X_{S'(1)}, \dots, X_{S'(K)}; Y)$$

CMIM, mRMR (and others)

Compromise Joint Informativeness for Tractability

$$\{X_1, \dots, X_D\} \in \mathbb{R}, Y \in \{1 \dots C\}$$

$$S = \operatorname{argmax}_{S', |S'|=K} I(X_{S'(1)}, X_{S'(2)}, \dots, X_{S'(K)}; Y)$$

Calculate  $I(X; Y)$  using a joint law

- ▶ Differential Entropy

$$h(X) = - \int_{\mathcal{X}} p(x) \log(p(x)) dx$$

- ▶ Relative Entropy (Kullback-Leibler Divergence )

$$D(p\|q) = \int_{\mathcal{X}} p(x) \log \left( \frac{p(x)}{q(x)} \right) dx$$

- ▶ Mutual Information

$$I(X; Y) = \int_{\mathcal{X}} \int_{\mathcal{Y}} p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right) dx dy$$

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= \underbrace{H(X)}_{\text{intractable}} - \sum_Y P(Y = y) \underbrace{H(X|Y = y)}_{\text{intractable}} \end{aligned}$$

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= \underbrace{H(X)}_{\text{intractable}} - \sum_Y P(Y = y) \underbrace{H(X|Y = y)}_{\text{intractable}} \end{aligned}$$

Parametric model  $p_{X|Y}$

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= \underbrace{H(X)}_{\text{intractable}} - \sum_Y P(Y = y) \underbrace{H(X|Y = y)}_{\text{intractable}} \end{aligned}$$

Parametric model  $p_{X|Y} = \mathcal{N}(\mu_y, \Sigma_y)$

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= \underbrace{H(X)}_{\text{intractable}} - \sum_Y P(Y = y) \underbrace{H(X|Y = y)}_{\text{tractable}} \end{aligned}$$

Parametric model  $p_{X|Y} = N(\mu_y, \Sigma_y)$

$$H(X|Y) = \frac{1}{2} \log(|\Sigma_y|) + \frac{n}{2} (\log 2\pi + 1).$$



Given

$$E(x) = \mu, E\left((x - \mu)(x - \mu)^T\right) = \Sigma$$

then the multivariate Normal

$$\vec{x} \sim N(\vec{\mu}, \Sigma)$$

is the Maximum Entropy Distribution

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= \underbrace{H(X)}_{\text{intractable}} - \sum_Y P(Y = y) \underbrace{H(X|Y = y)}_{\text{tractable}} \end{aligned}$$

Parametric model  $p_{X|Y} = \mathcal{N}(\mu_y, \Sigma_y)$

$$H(X|Y) = \frac{1}{2} \log(|\Sigma_y|) + \frac{n}{2} (\log 2\pi + 1).$$

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= \underbrace{H(X)}_{\text{intractable}} - \sum_Y P(Y = y) \underbrace{H(X|Y = y)}_{\text{tractable}} \end{aligned}$$

$$H(X) = H \left( \sum_y \pi_y N(\mu_y, \Sigma_y) \right)$$

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= \underbrace{H(X)}_{\text{intractable}} - \sum_Y P(Y = y) \underbrace{H(X|Y = y)}_{\text{tractable}} \end{aligned}$$

$$H(X) = H\left(\sum_y \pi_y N(\mu_y, \Sigma_y)\right) \rightarrow \underbrace{\text{Entropy of Mixture of Gaussians}}_{\text{No Analytical Solution}}$$

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= \underbrace{H(X)}_{\text{intractable}} - \sum_Y P(Y = y) \underbrace{H(X|Y = y)}_{\text{tractable}} \end{aligned}$$

$$H(X) = H\left(\sum_y \pi_y N(\mu_y, \Sigma_y)\right) \rightarrow \underbrace{\text{Entropy of Mixture of Gaussians}}_{\text{No Analytical Solution}}$$

Upper Bound or Approximate

$$H\left(\sum_y \pi_y N(\mu_y, \Sigma_y)\right)$$

## A Normal Upper Bound



$$p_X = \sum_Y \pi_Y p_{X|Y}$$

$$p_X = \sum_Y \pi_y p_{X|Y}$$

$$p^* \sim N(\mu^*, \Sigma^*) \xrightarrow{\text{maxEnt}} H(p_X) \leq H(p^*)$$

$$p_X = \sum_Y \pi_Y p_{X|Y}$$

$$p^* \sim N(\mu^*, \Sigma^*) \xrightarrow{\text{maxEnt}} H(p_X) \leq H(p^*)$$

$$I(X; Y) \leq H(p^*) - \sum_Y P_Y H(X|Y)$$



$$p_X = \sum_Y \pi_y p_{X|Y}$$

$$p^* \sim N(\mu^*, \Sigma^*) \xrightarrow{\max Ent} H(p_X) \leq H(p^*)$$

$$I(X; Y) \leq H(p^*) - \sum_Y P_Y H(X|Y)$$

$$I(X; Y) \leq H(Y) = - \sum_y p_y \log p_y$$

## A Normal Upper Bound



$$p_X = \sum_Y \pi_y p_{X|Y}$$

$$I(X; Y) \leq H(p^*) - \sum_Y P_Y H(X|Y)$$

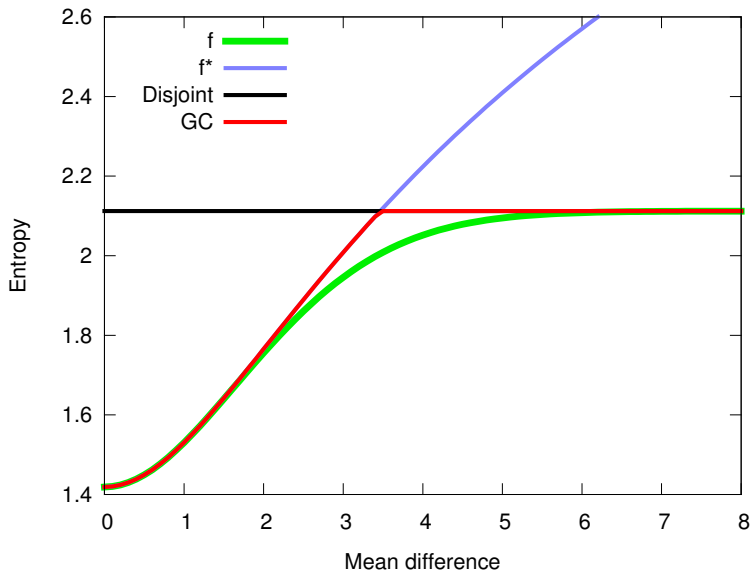
$$I(X; Y) \leq H(Y) = - \sum_y p_y \log p_y$$

$$p_X = \sum_Y \pi_y p_{X|Y}$$

$$I(X; Y) \leq H(p^*) - \sum_Y P_Y H(X|Y)$$

$$I(X; Y) \leq H(Y) = - \sum_y p_y \log p_y$$

$$I(X; Y) \leq \min \left( H(p^*), \sum_Y P_Y (H(X|Y) - \log(P_Y)) \right) - \sum_Y P_Y H(X|Y)$$



$$p_X = \sum_Y \pi_y p_{X|Y}$$

Under mild assumptions

$$\forall y, H(p^*) > H(p_{X|Y=y})$$

we can use an approximation to  $I(X; Y)$

$$\tilde{I}(X; Y) = \sum_y \min(H(p^*), H(X|Y) - \log p_y) p_y - \sum_y H(X|Y) p_y$$

## Mutual Information Approximation

$$S = \operatorname{argmax}_{S', |S'|=K} \left( \tilde{I}(X_{S'(1)}, X_{S'(2)}, \dots, X_{S'(K)}; Y) \right)$$

# Forward Selection



```
 $S_0 = \emptyset$   
for  $k = 1 \dots K$  do  
   $s^* = 0$   
  for  $X_j \in F \setminus S_{k-1}$  do  
     $S' \leftarrow S_{k-1} \cup X_j$   
     $z \leftarrow \hat{I}(S'; Y)$   
    if  $z > z^*$  then  
       $s^* \leftarrow s$   
       $S^* \leftarrow S'$   
    end if  
  end for  
   $S_i \leftarrow S^*$   
end for  
return  $S_K$ 
```

$$\hat{I}(S'; Y) \propto \sum_y \min(\log(|\Sigma^*|), \log(|\Sigma_y|) - \log p_y) p_y - \sum_y \log(|\Sigma_y|) p_y$$

At iteration  $k$  we need to calculate  $\forall j \in F \setminus S_{k-1}$

$$|\Sigma_{S_{k-1} \cup X_j}|$$



At iteration  $k$  we need to calculate  $\forall j \in F \setminus S_{k-1}$

$$|\Sigma_{S_{k-1} \cup X_j}|$$

The cost of calculating each determinant is  $O(k^3)$

```
 $S_0 = \emptyset$   
for  $k = 1 \dots K$  do  
   $s^* = 0$   
  for  $X_j \in F \setminus S_{k-1}$  do  
     $S' \leftarrow S_{k-1} \cup X_j$   
     $z \leftarrow \tilde{I}(S'; Y)$   
    if  $z > z^*$  then  
       $s^* \leftarrow s$   
       $S^* \leftarrow S'$   
    end if  
  end for  
   $S_i \leftarrow S^*$   
end for  
return  $S_K$ 
```

Overall Complexity  $O(|Y||F|K^4)$

```
 $S_0 = \emptyset$   
for  $k = 1 \dots K$  do  
   $s^* = 0$   
  for  $X_j \in F \setminus S_{k-1}$  do  
     $S' \leftarrow S_{k-1} \cup X_j$   
     $z \leftarrow \tilde{I}(S'; Y)$   
    if  $z > z^*$  then  
       $s^* \leftarrow s$   
       $S^* \leftarrow S'$   
    end if  
  end for  
   $S_i \leftarrow S^*$   
end for  
return  $S_K$ 
```

Overall Complexity  $O(|Y||F|K^4)$

However...

$$\Sigma_{S_{k-1} \cup X_j} = \begin{bmatrix} \Sigma_{S_{k-1}} & \Sigma_{j, S_{k-1}} \\ \Sigma_{j, S_{k-1}}^T & \sigma_j^2 \end{bmatrix}$$

$$\Sigma_{S_{k-1} \cup X_j} = \begin{bmatrix} \Sigma_{S_{k-1}} & \Sigma_{j, S_{k-1}} \\ \Sigma_{j, S_{k-1}}^T & \sigma_j^2 \end{bmatrix}$$

We can exploit the matrix determinant lemma (twice)

$$|\Sigma + uv^T| = (1 + v^T \Sigma^{-1} u) |\Sigma|$$

To compute each determinant in  $O(n^2)$

```
 $S_0 = \emptyset$   
for  $k = 1 \dots K$  do  
   $s^* = 0$   
  for  $X_j \in F \setminus S_{k-1}$  do  
     $S' \leftarrow S_{k-1} \cup X_j$   
     $z \leftarrow \tilde{I}(S'; Y)$   
    if  $z > z^*$  then  
       $s^* \leftarrow s$   
       $S^* \leftarrow S'$   
    end if  
  end for  
   $S_i \leftarrow S^*$   
end for  
return  $S_K$ 
```

Overall Complexity  $O(|Y||F|K^3)$

```
 $S_0 = \emptyset$   
for  $k = 1 \dots K$  do  
   $s^* = 0$   
  for  $X_j \in F \setminus S_{k-1}$  do  
     $S' \leftarrow S_{k-1} \cup X_j$   
     $z \leftarrow \tilde{I}(S'; Y)$   
    if  $z > z^*$  then  
       $s^* \leftarrow s$   
       $S^* \leftarrow S'$   
    end if  
  end for  
   $S_i \leftarrow S^*$   
end for  
return  $S_K$ 
```

Overall Complexity  $O(|Y||F|K^3)$

Faster?

$$I(X, Z; Y) = I(Z; Y) + I(X; Y | Z)$$

$$I(S'_k; Y) = I(X_j; Y | S_{k-1}) + \underbrace{I(S_{k-1}; Y)}_{\text{common}}$$



$$I(X, Z; Y) = I(Z; Y) + I(X; Y | Z)$$

$$I(S'_k; Y) = I(X_j; Y | S_{k-1}) + \underbrace{I(\cancel{S_{k-1}}; Y)}_{\text{common}}$$

$$\operatorname{argmax}_{X_j \in F \setminus S_{k-1}} I(X_j; Y | S_{k-1}) =$$

$$\operatorname{argmax}_{X_j \in F \setminus S_{k-1}} (H(X_j | S_{k-1}) - H(X_j | Y, S_{k-1}))$$

$$\begin{aligned} H(X_j | S_{k-1}) &= \int_{\mathbb{R}^{|S_{k-1}|}} H(X_j | S_{k-1} = s) \mu_{S_{k-1}}(s) ds \\ &= \frac{1}{2} \log \sigma_{j|S_{k-1}}^2 + \frac{1}{2} (\log 2\pi + 1) \end{aligned}$$

$$\sigma_{j|S_{k-1}}^2 = \sigma_j^2 - \Sigma_{j,S_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{j,S_{k-1}}.$$

$$\sigma_j^2|S_{k-1} = \sigma_j^2 - \Sigma_{j,S_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{j,S_{k-1}}$$

Assume  $X_j$  was chosen at iteration  $k - 1$

$$\begin{aligned} \Sigma_{S_{k-1}} &= \begin{bmatrix} \Sigma_{S_{k-2}} & \Sigma_{i,S_{k-2}} \\ \Sigma_{i,S_{k-2}}^T & \sigma_i^2 \end{bmatrix} \\ &= \begin{bmatrix} \Sigma_{S_{k-2}} & 0_{k-2} \\ 0_{n-2}^T & \sigma_i^2 \end{bmatrix} + \mathbf{e}_{n+1} \Sigma_{i,S_{k-2}}^T + \Sigma_{i,S_{k-2}} \mathbf{e}_{n+1}^T \end{aligned}$$

From Sherman-Morrison formula

$$\left(\Sigma + uv^T\right)^{-1} = \Sigma^{-1} - \frac{\Sigma^{-1}uv^T\Sigma^{-1}}{1 + v^T\Sigma^{-1}u}$$

From Sherman-Morrison formula

$$\left(\Sigma + uv^T\right)^{-1} = \Sigma^{-1} - \frac{\Sigma^{-1}uv^T\Sigma^{-1}}{1 + v^T\Sigma^{-1}u}$$

$$\Sigma_{S_{n-1}}^{-1} = \begin{bmatrix} \Sigma_{S_{n-2}}^{-1} & -\frac{1}{\beta\sigma_i^2}u \\ -\frac{1}{\beta\sigma_i^2}u^T & \frac{1}{\beta\sigma_i^2} \end{bmatrix} + \frac{1}{\beta\sigma_i^2} \begin{bmatrix} u \\ 0 \end{bmatrix} \begin{bmatrix} u^T & 0 \end{bmatrix}$$

$$\Sigma_{S_{n-1}}^{-1} = \begin{bmatrix} \Sigma_{S_{n-2}}^{-1} & -\frac{1}{\beta\sigma_i^2} \mathbf{u} \\ -\frac{1}{\beta\sigma_i^2} \mathbf{u}^T & \frac{1}{\beta\sigma_i^2} \end{bmatrix} + \frac{1}{\beta\sigma_i^2} \begin{bmatrix} \mathbf{u} \\ 0 \end{bmatrix} \begin{bmatrix} \mathbf{u}^T & 0 \end{bmatrix}$$

$$\begin{aligned} \sigma_{j|S_{n-1}}^2 &= \sigma_j^2 - \overbrace{\Sigma_{j,S_{n-2}}^T \Sigma_{S_{n-2}}^{-1} \Sigma_{j,S_{n-2}}}^{\text{Previous Round}} \in O(n^2) \\ &+ \frac{\sigma_{ji}^2}{\beta\sigma_i^2} \mathbf{u}^T \Sigma_{j,S_{n-2}} - \Sigma_{j,S_{n-1}}^T \begin{bmatrix} -\frac{1}{\beta\sigma_i^2} \mathbf{u} \\ \frac{1}{\beta\sigma_i^2} \end{bmatrix} \sigma_{ji}^2 \in O(n) \\ &- \frac{1}{\beta\sigma_i^2} \left( \Sigma_{j,S_{n-1}}^T \begin{bmatrix} \mathbf{u} \\ 0 \end{bmatrix} \right) \left( \begin{bmatrix} \mathbf{u}^T & 0 \end{bmatrix} \Sigma_{j,S_{n-1}} \right) \in O(n) \end{aligned}$$

```
for  $k = 1 \dots K$  do
   $s^* = 0$ 
  for  $X_j \in F \setminus S_{k-1}$  do
     $S' \leftarrow S_{k-1} \cup X_j$ 
     $z \leftarrow \hat{I}(S')$ 
    if  $z > z^*$  then
       $s^* \leftarrow z$ 
       $S^* \leftarrow S'$ 
    end if
  end for
   $S_k \leftarrow S^*$ 
end for
return  $S_K$ 
```

Overall Complexity  $O(|Y||F|K^2)$

Faster!



```
for  $k = 1 \dots K$  do
   $s^* = 0$ 
  for  $X_j \in F \setminus S_{k-1}$  do
     $S' \leftarrow S_{k-1} \cup X_j$ 
     $z \leftarrow \hat{I}(S')$ 
    if  $z > z^*$  then
       $s^* \leftarrow s$ 
       $S^* \leftarrow S'$ 
    end if
  end for
   $S_i \leftarrow S^*$ 
end for
return  $S_K$ 
```

Overall Complexity  $O(|Y||F|K^2)$

Even Faster?



Main bottleneck  $O(k)$

$$\sigma_{j|S_{k-1}}^2 = \sigma_j^2 - \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}}$$

Main bottleneck  $O(k)$

$$\sigma_{j|S_{k-1}}^2 = \sigma_j^2 - \sum_{j \in S_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \sum_{j \in S_{k-1}}$$

Skip non-promising features

```
 $S_0 = \emptyset$   
for  $k = 1 \dots K$  do  
   $s^* = 0$   
  for  $X_j \in F \setminus S_{k-1}$  do  
     $S' \leftarrow S_{k-1} \cup X_j$   
     $z \leftarrow \hat{I}(S')$   
    if  $z > z^*$  then  
       $s^* \leftarrow s$   
       $S^* \leftarrow S'$   
    end if  
  end for  
   $S_i \leftarrow S^*$   
end for  
return  $S_K$ 
```

Cheap ( $O(1)$ ) score  $c$ :

if  $c < z^*$  then  $z < z^*$

...

$z^* = 0$

**for**  $X_j \in F \setminus S_{k-1}$  **do**

$S' \leftarrow S_{k-1} \cup X_j$

$c \leftarrow ?$

**if**  $c > z^*$  **then**

$z \leftarrow \hat{I}(S')$

**if**  $z > z^*$  **then**

$s^* \leftarrow s$

$S^* \leftarrow S'$

**end if**

**end if**

**end for**

...

Cheap ( $O(1)$ ) score  $c$ :

if  $c < z^*$  then  $z < z^*$

$$\sigma_{j|S_{k-1}}^2 = \sigma_j^2 - \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}}$$

$$\Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}} = \Sigma_{jS_{k-1}}^T U \Lambda U^T \Sigma_{jS_{k-1}}$$

$$\|\Sigma_{jS_{k-1}}^T\|_2^2 \max_i \lambda_i \geq \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}} \geq \underbrace{\|\Sigma_{jS_{k-1}}^T\|_2^2 \min_i \lambda_i}_{O(1)}$$

## An $O(1)$ Bound



$$\sigma_{j|S_{k-1}}^2 = \sigma_j^2 - \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}}$$

$$\Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}} = \Sigma_{jS_{k-1}}^T U \Lambda U^T \Sigma_{jS_{k-1}}$$

$$\|\Sigma_{jS_{k-1}}^T\|_2^2 \max_i \lambda_i \geq \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}} \geq \underbrace{\|\Sigma_{jS_{k-1}}^T\|_2^2 \min_i \lambda_i}_{O(1)}$$

$$\sigma_{j|S_{k-1}}^2 = \sigma_j^2 - \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}}$$

$$\Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}} = \Sigma_{jS_{k-1}}^T U \Lambda U^T \Sigma_{jS_{k-1}}$$

$$\|\Sigma_{jS_{k-1}}^T\|_2^2 \max_i \lambda_i \geq \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}} \geq \underbrace{\|\Sigma_{jS_{k-1}}^T\|_2^2 \min_i \lambda_i}_{O(1)}$$

$$\sigma_{j|S_{k-1}}^2 = \sigma_j^2 - \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}}$$

$$\Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}} = \Sigma_{jS_{k-1}}^T U \Lambda U^T \Sigma_{jS_{k-1}}$$

$$\|\Sigma_{jS_{k-1}}^T\|_2^2 \max_i \lambda_i \geq \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}} \geq \underbrace{\|\Sigma_{jS_{k-1}}^T\|_2^2 \min_i \lambda_i}_{O(1)}$$



$$\sigma_{j|S_{k-1}}^2 = \sigma_j^2 - \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}}$$

$$\Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}} = \Sigma_{jS_{k-1}}^T U \Lambda U^T \Sigma_{jS_{k-1}}$$

$$\|\Sigma_{jS_{k-1}}^T\|_2^2 \max_i \lambda_i \geq \Sigma_{jS_{k-1}}^T \Sigma_{S_{k-1}}^{-1} \Sigma_{jS_{k-1}} \geq \underbrace{\|\Sigma_{jS_{k-1}}^T\|_2^2 \min_i \lambda_i}_{O(1)}$$

# EigenSystem Update Problem



Given  $U^n, \Lambda^n$  such that

$$\Sigma^n = U^{nT} \Lambda^n U^n$$

# EigenSystem Update Problem



Given  $U^n, \Lambda^n$  such that

$$\Sigma^n = U^{nT} \Lambda^n U^n$$

Find  $U^{n+1}, \Lambda^{n+1}$

$$\Sigma^{n+1} = \begin{bmatrix} \Sigma^n & \mathbf{v} \\ \mathbf{v}^T & 1 \end{bmatrix} = U^{n+1T} \Lambda^{n+1} U^{n+1}$$

assume  $\Sigma^{n+1} \in S_{++}^{n+1}$

$$U^n = [ \mathbf{u}_1^n \quad \mathbf{u}_2^n \quad \dots \quad \mathbf{u}_n^n ], \quad \Lambda^n = \begin{bmatrix} \lambda_1^n & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_n^n \end{bmatrix}$$

$$U^n = [ \mathbf{u}_1^n \quad \mathbf{u}_2^n \quad \dots \quad \mathbf{u}_n^n ], \quad \Lambda^n = \begin{bmatrix} \lambda_1^n & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_n^n \end{bmatrix}$$

$$\begin{bmatrix} \Sigma^n & \mathbf{0} \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{u}_i^n \\ 0 \end{bmatrix} = \lambda_i^n \underbrace{\begin{bmatrix} \mathbf{u}_i^n \\ 0 \end{bmatrix}}_{\mathbf{u}_i'^n}$$

$$U^n = [ \mathbf{u}_1^n \quad \mathbf{u}_2^n \quad \dots \quad \mathbf{u}_n^n ], \quad \Lambda^n = \begin{bmatrix} \lambda_1^n & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_n^n \end{bmatrix}$$

$$\begin{bmatrix} \Sigma^n & \mathbf{0} \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{u}_i^n \\ 0 \end{bmatrix} = \lambda_i^n \underbrace{\begin{bmatrix} \mathbf{u}_i^n \\ 0 \end{bmatrix}}_{\mathbf{u}_i'^n}$$

$$\begin{bmatrix} \Sigma^n & \mathbf{0} \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ 1 \end{bmatrix}$$

$$\underbrace{\begin{bmatrix} \Sigma^n & \mathbf{0} \\ \mathbf{0}^T & 1 \end{bmatrix}}_{\Sigma'} = \begin{bmatrix} \mathbf{u}'_1{}^T \\ \vdots \\ \mathbf{e}^{n+1}{}^T \end{bmatrix} \underbrace{\begin{bmatrix} \lambda_1^n & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{bmatrix}}_{\Lambda'} \underbrace{\begin{bmatrix} \mathbf{u}'_1{}^n & \cdots & \mathbf{e}^{n+1} \end{bmatrix}}_{U'}$$

$$\Sigma^{n+1} = \Sigma' + \mathbf{e}_{n+1} \begin{bmatrix} \mathbf{v} \\ 0 \end{bmatrix}^T + \begin{bmatrix} \mathbf{v} \\ 0 \end{bmatrix} \mathbf{e}_{n+1}^T$$

$$\left( \Sigma' + \mathbf{e}_{n+1} \begin{bmatrix} \mathbf{v} \\ 0 \end{bmatrix}^T + \begin{bmatrix} \mathbf{v} \\ 0 \end{bmatrix} \mathbf{e}_{n+1}^T \right) \mathbf{u}^{n+1} = \lambda^{n+1} \mathbf{u}^{n+1}$$

$$\left( \Sigma' + \mathbf{e}_{n+1} \begin{bmatrix} \mathbf{v} \\ 0 \end{bmatrix}^T + \begin{bmatrix} \mathbf{v} \\ 0 \end{bmatrix} \mathbf{e}_{n+1}^T \right) u^{n+1} = \lambda^{n+1} u^{n+1}$$

$$U'^T \left( \Sigma' + \mathbf{e}_{n+1} \begin{bmatrix} \mathbf{v} \\ 0 \end{bmatrix}^T + \begin{bmatrix} \mathbf{v} \\ 0 \end{bmatrix} \mathbf{e}_{n+1}^T \right) \underbrace{U' U'^T}_{U' U'^T = I} u^{n+1} = \lambda^{n+1} U'^T u^{n+1}$$

$$\stackrel{((4))}{\implies} \left( \Lambda' + \mathbf{e}_{n+1} \mathbf{q}^T + \mathbf{q} \mathbf{e}_{n+1}^T \right) U'^T u^{n+1} = \lambda^{n+1} U'^T u^{n+1}$$

$$U'^T \Sigma' U' = \Lambda' \tag{4}$$



$$\underbrace{\left( \Lambda' + \mathbf{e}_{n+1} \mathbf{q}^T + \mathbf{q} \mathbf{e}_{n+1}^T \right)}_{\Sigma''} U'^T u^{n+1} = \lambda^{n+1} U'^T u^{n+1}$$

→  $\Sigma^{n+1}$  and  $\Sigma''$  share eigenvalues.

→  $U^{n+1} = U' U''$

$$|\Sigma'' - \lambda I| = \prod_j (\lambda'_j - \lambda) + \sum_{i < n+1} -q_i^2 \prod_{j \neq i, j < n+1} (\lambda'_j - \lambda)$$

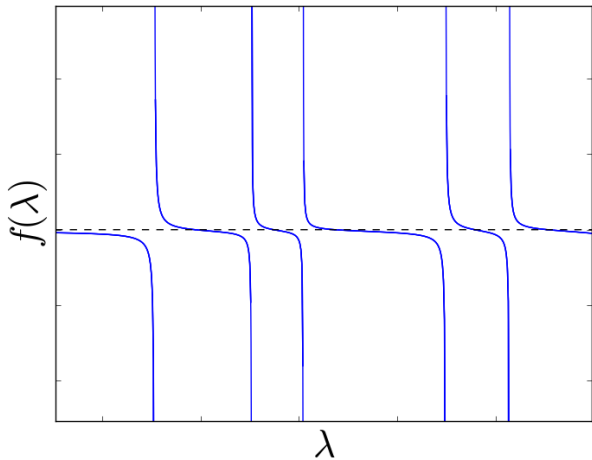
$$f(\lambda) = \lambda'_{n+1} - \lambda + \sum_i \frac{-q_i^2}{(\lambda'_i - \lambda)}$$

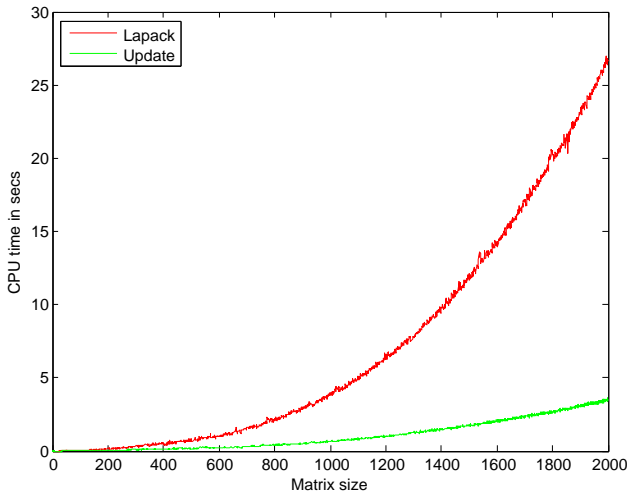
$$|\Sigma'' - \lambda I| = \prod_j (\lambda'_j - \lambda) + \sum_{i < n+1} -q_i^2 \prod_{j \neq i, j < n+1} (\lambda'_j - \lambda)$$

$$f(\lambda) = \lambda'_{n+1} - \lambda + \sum_i \frac{-q_i^2}{(\lambda'_i - \lambda)}$$

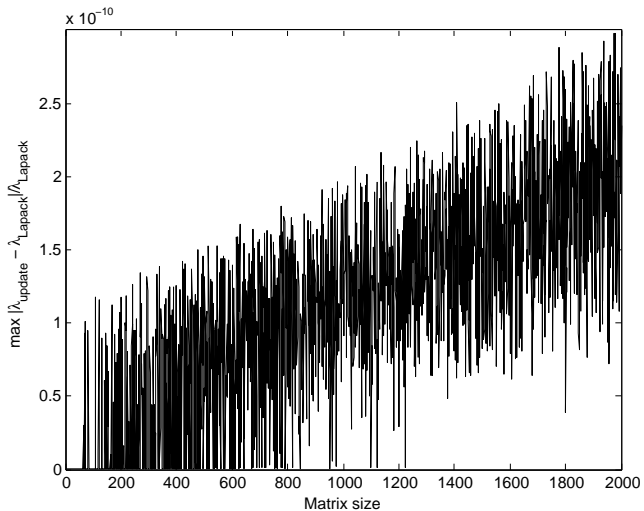
$$\forall i, \lim_{\lambda \xrightarrow{\geq} \lambda'_i} f(\lambda) = +\infty, \lim_{\lambda \xrightarrow{\leq} \lambda'_i} f(\lambda) = -\infty$$

$$\frac{\partial f(\lambda)}{\partial \lambda} = -1 + \sum_i \frac{-q_i^2}{(\lambda'_i - \lambda)^2} \leq 0$$





Comparison between scratch and update



Numerical stability (eigenvalues)

Set to go ...



Now that all the machinery is in place, How did we do?

## Nice Results

SVMLin	CIFAR			STL			INRIA		
	10	50	100	10	50	100	10	50	100
Fisher	25.19	39.47	48.12	26.09	34.63	38.02	92.55	94.03	94.68
FCBF	<b>33.65</b>	47.77	54.97	31.74	38.11	40.66	<b>94.14</b>	<b>96.03</b>	96.03
MRMR	27.94	37.78	43.63	28.26	31.16	33.12	86.03	86.77	86.72
SBMLR	30.43	<b>51.41</b>	<b>56.81</b>	<b>32.29</b>	43.29	47.15	85.92	88.57	88.64
tTest	25.69	40.17	45.12	26.72	36.23	39.14	80.01	87.64	89.23
InfoGain	24.79	37.98	47.37	27.17	33.70	37.84	92.35	93.75	94.68
Spec. Clus.	17.19	32.78	42.6	18.91	32.65	38.24	<b>92.67</b>	93.64	94.44
RelieFF	24.56	38.17	46.51	29.16	38.05	42.94	90.99	<b>95.97</b>	<b>96.36</b>
CFS	31.49	42.17	51.70	28.63	38.54	41.88	88.64	96.11	<b>97.53</b>
CMTF	21.10	40.39	47.71	27.61	38.99	42.32	79.09	89.49	93.01
<i>BAHSIC</i>	-	-	-	28.95	39.05	45.49	78.54	89.77	91.96
<b>GC.E</b>	32.45	50.15	55.06	31.20	43.31	<b>49.75</b>	87.73	91.96	93.13
<b>GC. MI</b>	<b>36.47</b>	<b>51.44</b>	<b>55.39</b>	<b>32.50</b>	<b>44.15</b>	<b>48.88</b>	89.76	<b>95.71</b>	<b>96.45</b>
<b>GKL.E</b>	<b>37.51</b>	<b>52.11</b>	<b>56.41</b>	<b>33.44</b>	<b>44.27</b>	<b>50.54</b>	85.31	92.05	<b>96.36</b>
<b>GKL. MI</b>	33.71	47.17	51.12	32.16	<b>44.87</b>	47.96	85.66	92.14	95.16

GC.MI was the fastest of the more complex algorithms



We use estimates

$$\hat{\Sigma}_N = \frac{1}{N} P^T P$$

For (sub)-Gaussian data we have<sup>2</sup>

$$\text{If } N \geq C(t/\epsilon)^2 d \quad \text{then } \|\hat{\Sigma}_N - \Sigma\| \leq \epsilon$$

---

<sup>2</sup>with probability at least  $1 - 2e^{-ct^2}$

We use estimates

$$\hat{\Sigma}_N = \frac{1}{N} P^T P$$

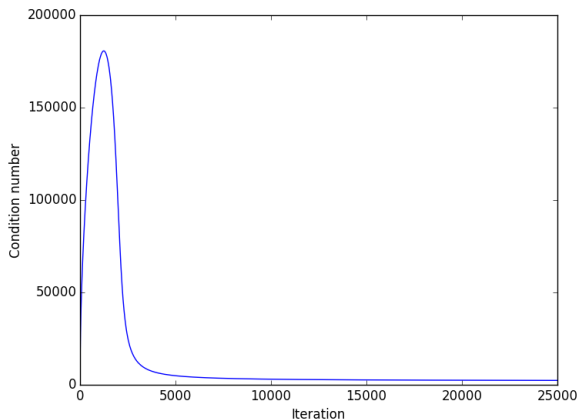
For (sub)-Gaussian data we have<sup>2</sup>

$$\text{If } N \geq C(t/\epsilon)^2 d \quad \text{then } \|\hat{\Sigma}_N - \Sigma\| \leq \epsilon$$

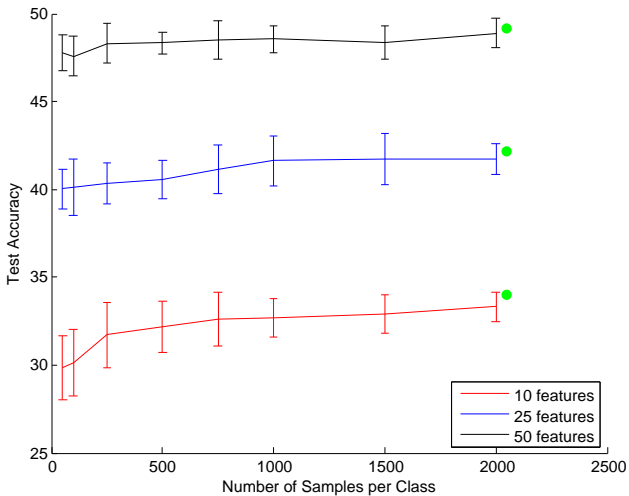
However the faster implementations use  $\hat{\Sigma}_N^{-1}$

---

<sup>2</sup>with probability at least  $1 - 2e^{-ct^2}$



$$\kappa(\Sigma) = \frac{\frac{\|\Sigma^{-1}e\|}{\|\Sigma^{-1}b\|}}{\frac{\|e\|}{\|b\|}} = \|\Sigma^{-1}\| \|\Sigma\|, \text{ for } d = 2048 \text{ and various values of } N$$



Effect of sample size on performance when using the Gaussian Approximation for the CIFAR dataset.

# Jointly Informative Feature Selection Made Tractable by Gaussian Modeling

L.Lefakis and F.Fleuret

Journal of Machine Learning Research, 2016

The End



Thank You!