

# CS480/680: Introduction to Machine Learning

## Lec 02: Linear Regression

Yaoliang Yu



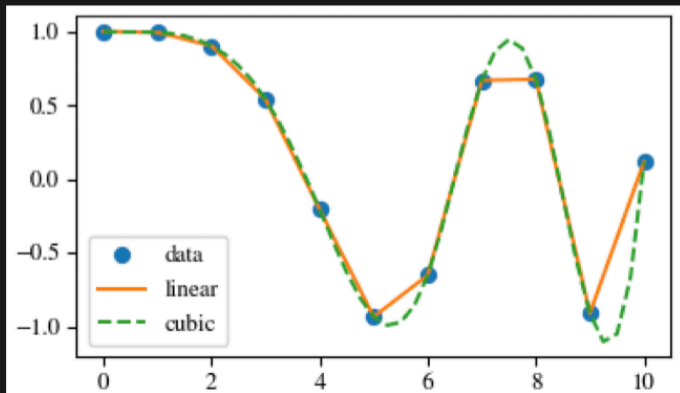
UNIVERSITY OF  
**WATERLOO**

FACULTY OF MATHEMATICS  
**DAVID R. CHERITON SCHOOL  
OF COMPUTER SCIENCE**

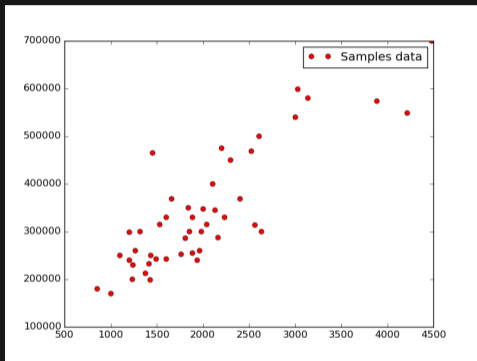
Jan 14, 2025

# Regression

- Given training data  $\{(\mathbf{x}_i, \mathbf{y}_i)\}$ , find  $f: \mathcal{X} \rightarrow \mathcal{Y}$  such that  $f(\mathbf{x}_i) \approx y_i$ 
  - $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^d$ : feature vector for the  $i$ -th training example
  - $\mathbf{y}_i \in \mathcal{Y} \subseteq \mathbb{R}^t$ :  $t$  responses, e.g.  $t = 1$  or even  $t = \infty$

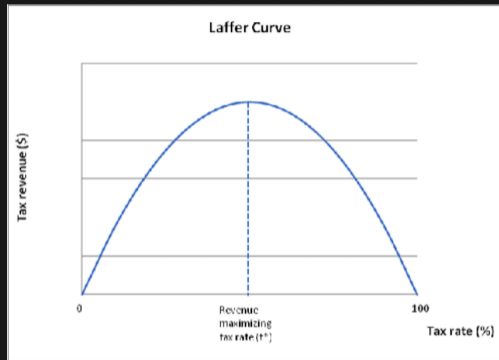
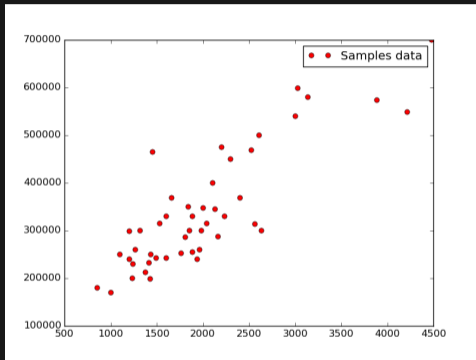


# Some Examples



- Prior knowledge on the functional form of  $f$
- Linear vs. nonlinear

# Some Examples



- Prior knowledge on the functional form of  $f$
- Linear vs. nonlinear

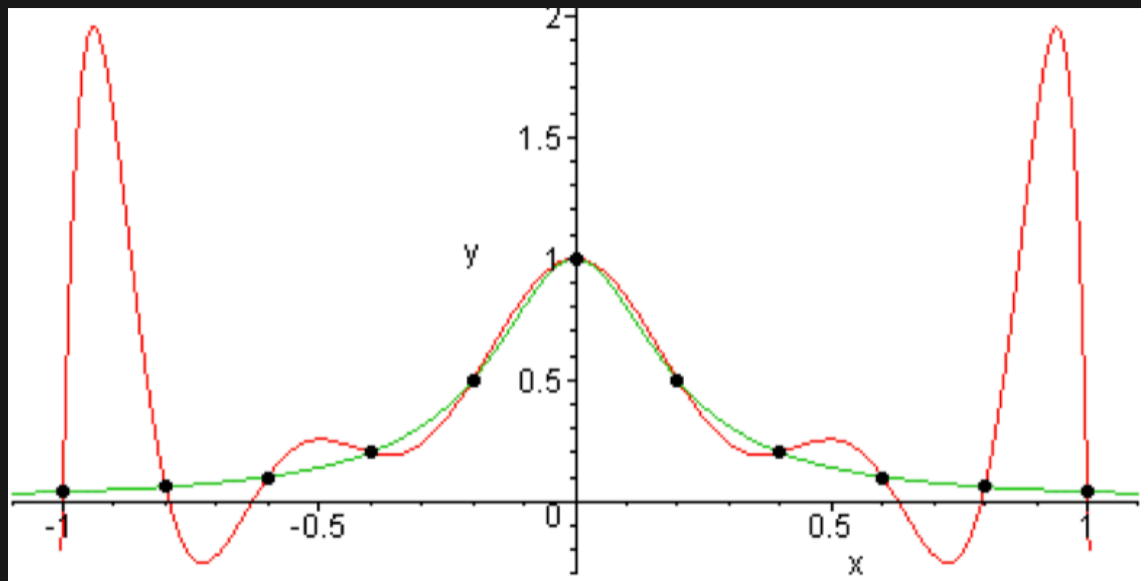
# The Difficulty

Theorem: Exact interpolation is always possible

For **any**\* **finite** training data  $\{(x_i, y_i) : i = 1, \dots, n\}$ , there exist **infinitely** many functions  $f$  such that for all  $i$ ,

$$f(x_i) = y_i.$$

- No amount of training data is enough to decide on a unique  $f$ !
- On new data  $x$ , our prediction  $\hat{y} = f(x)$  can vary wildly!
- This is where prior knowledge of  $f$  comes into play
- **Occam's razor**: “the simplest explanation is usually the correct one”

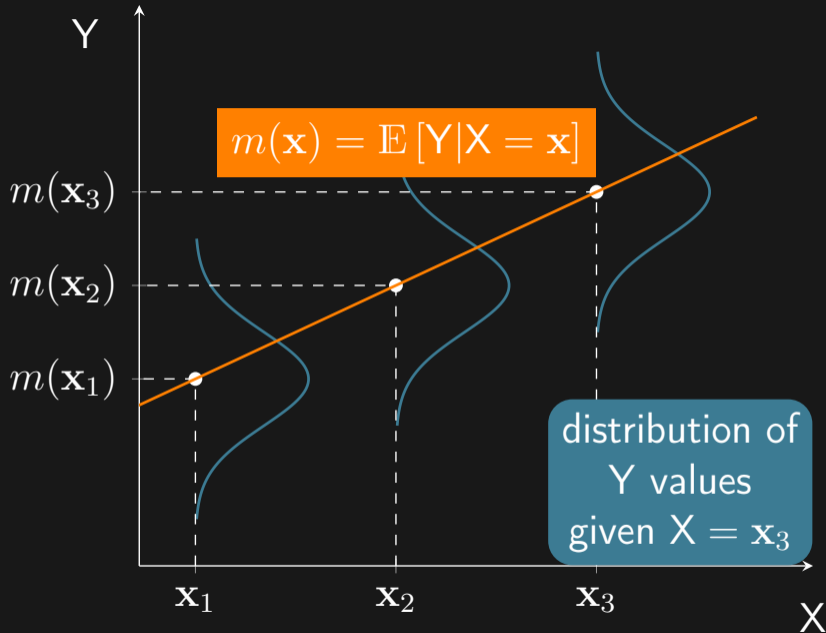


# Statistical Learning

- Training and test data are both iid samples from the **same unknown** distribution  $\mathbb{P}$

$$- (X_i, Y_i) \sim \mathbb{P} \text{ and } (X, Y) \sim \mathbb{P}$$

- **Least squares** regression:  $\min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E} \|f(\mathbf{X}) - Y\|_2^2$
- **Regression function**:  $m(\mathbf{x}) = \mathbb{E}[Y | X = \mathbf{x}]$
- Need to know the distribution  $\mathbb{P}$ , i.e., **all** pairs  $(X, Y)$ !
- Changing the square loss changes the regression function accordingly





# Bias-Variance Decomposition

$$\begin{aligned}\mathbb{E}\|f(\mathbf{X}) - \mathbf{Y}\|_2^2 &= \mathbb{E}\|f(\mathbf{X}) - m(\mathbf{X}) + m(\mathbf{X}) - \mathbf{Y}\|_2^2 \\ &= \mathbb{E}\|f(\mathbf{X}) - m(\mathbf{X})\|_2^2 + \mathbb{E}\|m(\mathbf{X}) - \mathbf{Y}\|_2^2 \\ &\quad + 2\mathbb{E}\langle f(\mathbf{X}) - m(\mathbf{X}), m(\mathbf{X}) - \mathbf{Y} \rangle \\ &= \underbrace{\mathbb{E}\|f(\mathbf{X}) - m(\mathbf{X})\|_2^2}_{\text{bias}^2} + \underbrace{\mathbb{E}\|m(\mathbf{X}) - \mathbf{Y}\|_2^2}_{\text{noise variance}}\end{aligned}$$

- The noise variance does not depend on our choice of  $f$ !
  - it is an inherent measure of the difficulty of our problem
- We aim to choose  $f \approx m$  to minimize bias hence squared error

# Sampling $\rightarrow$ Training

$$\min_{f:\mathcal{X}\rightarrow\mathcal{Y}} \hat{\mathbb{E}}\|f(\mathbf{X}) - \mathbf{Y}\|_2^2 = \frac{1}{n} \sum_{i=1}^n \|f(\mathbf{X}_i) - \mathbf{Y}_i\|_2^2$$

- Replace expectation with sample average:  $(\mathbf{X}_i, \mathbf{Y}_i) \stackrel{i.i.d.}{\sim} P$
- Finite training set  $\rightarrow$  exact interpolation paradox!
- Need to restrict the form of  $f$ , using prior knowledge
- (Uniform) law of large numbers: as training data size  $n \rightarrow \infty$ ,

$$\hat{\mathbb{E}} \rightarrow \mathbb{E} \text{ and (hopefully) } \operatorname{argmin} \hat{\mathbb{E}} \rightarrow \operatorname{argmin} \mathbb{E}$$

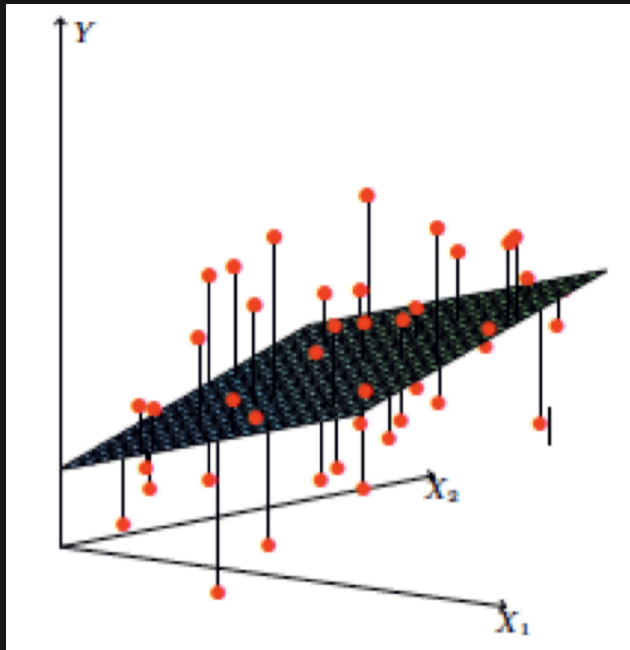
# Linear Least Squares Regression

- Affine function:  $f(\mathbf{x}) = W\mathbf{x} + \mathbf{b}$  with  $W \in \mathbb{R}^{t \times d}$  and  $\mathbf{b} \in \mathbb{R}^t$
- **Padding**:  $\mathbf{x} \leftarrow \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}$ ,  $\mathbf{W} \leftarrow [W, \mathbf{b}]$ , hence  $f(\mathbf{x}) = \mathbf{W}\mathbf{x}$
- In matrix form:  $\frac{1}{n} \sum_i \|f(\mathbf{x}_i) - \mathbf{y}_i\|_2^2 = \frac{1}{n} \|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_F^2$

–  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{(d+1) \times n}$ ,  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \mathbb{R}^{t \times n}$

–  $\|A\|_F = \sqrt{\sum_{ij} a_{ij}^2}$

$$\min_{\mathbf{W} \in \mathbb{R}^{t \times (d+1)}} \frac{1}{n} \|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_F^2$$



# Calculus Detour

- Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a smooth real-valued function
- Fix an inner product  $\langle \cdot, \cdot \rangle$
- Define the gradient  $\nabla f : \mathbb{R}^p \rightarrow \mathbb{R}^p$  as

$$\frac{df(\mathbf{w} + t\mathbf{z})}{dt} \Big|_{t=0} = \langle \nabla f(\mathbf{w}), \mathbf{z} \rangle$$

- LHS is the usual (scalar) derivative of the univariate function  $t \mapsto f(\mathbf{w} + t\mathbf{z})$
  - $\mathbf{w}$  and  $\mathbf{z}$  are fixed as constants: **directional derivative**
  - gradient  $\nabla f$  is **representation** of directional derivative over a chosen inner product
- **Chain rule** still holds

## Example: Univariate functions

Consider  $f : \mathbb{R} \rightarrow \mathbb{R}$  (i.e.,  $p = 1$ ) and the standard inner product  $\langle a, b \rangle := ab$ . By chain rule:

$$\frac{df(w + tz)}{dt} \Big|_{t=0} = f'(w + tz)z \Big|_{t=0} = f'(w)z = \langle f'(w), z \rangle,$$

i.e.,  $\nabla f(w) = f'(w)$ . What is the gradient if we choose  $\langle a, b \rangle := 2ab$ ?

## Example: Partial derivatives

Consider  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  and the standard inner product  $\langle \mathbf{w}, \mathbf{x} \rangle := \sum_j w_j x_j$ . Choose the direction  $\mathbf{z} = \mathbf{e}_j$  (i.e., 1 at the  $j$ -th entry and 0 elsewhere):

$$\frac{df(\mathbf{w} + t\mathbf{e}_j)}{dt} \Big|_{t=0} = \partial_j f(\mathbf{w}) = \langle \nabla f(\mathbf{w}), \mathbf{e}_j \rangle = [\nabla f(\mathbf{w})]_j,$$

i.e.,  $\nabla f(w) = [\partial_1 f(\mathbf{w}), \dots, \partial_p f(\mathbf{w})]$ .

## Example: Quadratic function

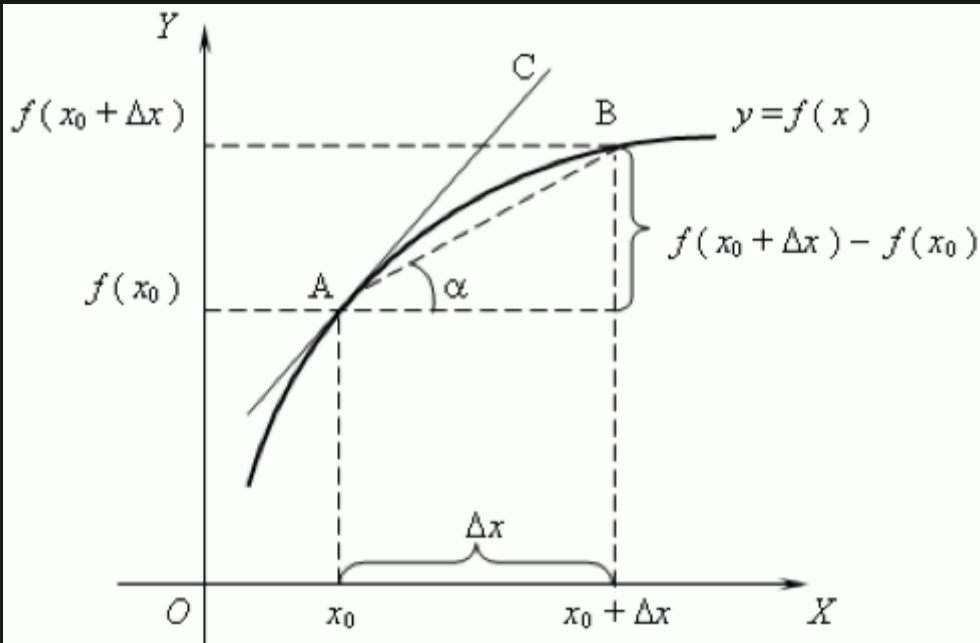
Consider the quadratic function  $f(\mathbf{w}) = \langle \mathbf{w}, A\mathbf{w} + \mathbf{b} \rangle + c$ .

$$\begin{aligned} f(\mathbf{w} + t\mathbf{z}) &= \langle \mathbf{w} + t\mathbf{z}, A(\mathbf{w} + t\mathbf{z}) + \mathbf{b} \rangle + c \\ &= t^2 \langle \mathbf{z}, A\mathbf{z} \rangle + t \langle \mathbf{w}, A\mathbf{z} \rangle + t \langle \mathbf{z}, A\mathbf{w} + \mathbf{b} \rangle + \langle \mathbf{w}, A\mathbf{w} + \mathbf{b} \rangle + c \end{aligned}$$

$$\frac{df(\mathbf{w} + t\mathbf{z})}{dt} \Big|_{t=0} = \langle \mathbf{w}, A\mathbf{z} \rangle + \langle \mathbf{z}, A\mathbf{w} + \mathbf{b} \rangle = \langle A^\top \mathbf{w} + A\mathbf{w} + \mathbf{b}, \mathbf{z} \rangle,$$

i.e.,  $\boxed{\nabla f(\mathbf{w}) = (A^\top + A)\mathbf{w} + \mathbf{b}}$ .

- $\langle \mathbf{a} + \mathbf{b}, \mathbf{x} + \mathbf{y} \rangle = \langle \mathbf{a}, \mathbf{x} \rangle + \langle \mathbf{a}, \mathbf{y} \rangle + \langle \mathbf{b}, \mathbf{x} \rangle + \langle \mathbf{b}, \mathbf{y} \rangle$
- $\langle \mathbf{a}, t\mathbf{b} \rangle = \langle t\mathbf{a}, \mathbf{b} \rangle = t \langle \mathbf{a}, \mathbf{b} \rangle$
- $\langle \mathbf{w}, A\mathbf{z} \rangle = \langle A^\top \mathbf{w}, \mathbf{z} \rangle, \langle A\mathbf{w}, \mathbf{z} \rangle = \langle \mathbf{w}, A^\top \mathbf{z} \rangle$

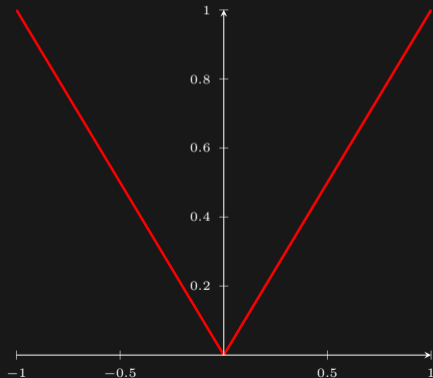
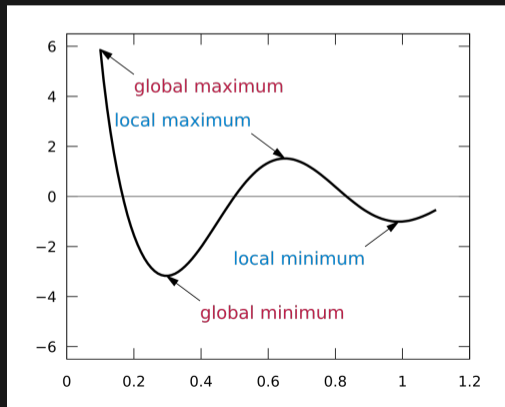




# Optimality Condition

Theorem: Fermat's necessary condition for extremity

If  $\mathbf{w}$  is a minimizer (or maximizer) of a differentiable function  $f$  over an open set, then  $f'(\mathbf{w}) = \mathbf{0}$ .



# Solving Linear Regression

$$\begin{aligned}\|\mathbf{WX} - \mathbf{Y}\|_F^2 &= \langle \mathbf{WX} - \mathbf{Y}, \mathbf{WX} - \mathbf{Y} \rangle \\ &= \langle \mathbf{W}, \mathbf{WXX}^\top - 2\mathbf{YX}^\top \rangle + \langle \mathbf{Y}, \mathbf{Y} \rangle\end{aligned}$$

- Taking derivative w.r.t.  $\mathbf{W}$  and setting to zero:

Normal equation  $\boxed{\mathbf{WXX}^\top = \mathbf{YX}^\top} \implies \mathbf{W} = \mathbf{YX}^\top (\mathbf{XX}^\top)^{-1} =: \mathbf{YX}^\dagger$

- $\mathbf{X} \in \mathbb{R}^{(d+1) \times n}$  hence  $\mathbf{XX}^\top \in \mathbb{R}^{(d+1) \times (d+1)}$ : may not be invertible if  $n \leq d + 1$ , but a solution always exists
- Even when invertible, **never compute the inverse directly!**
- Instead, solve the linear system or apply iterative gradient algorithm

# Prediction

- Once solved  $\mathbf{W}$  on the training set  $(\mathbf{X}, \mathbf{Y})$ , can predict on unseen data  $\mathbf{X}_{\text{test}}$ :

$$\hat{\mathbf{Y}}_{\text{test}} = \mathbf{W}\mathbf{X}_{\text{test}}$$

- We may evaluate our test error if true labels were available:

$$\frac{1}{n_{\text{test}}} \|\mathbf{Y}_{\text{test}} - \hat{\mathbf{Y}}_{\text{test}}\|_{\text{F}}^2$$

- We may compare to the training error:

$$\frac{1}{n} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_{\text{F}}^2, \quad \text{where } \hat{\mathbf{Y}} := \mathbf{W}\mathbf{X}$$

- Minimizing the training error as a means to reduce the test error
- Sometimes we even evaluate the test error using a different loss  $\mathbb{L}(\mathbf{Y}_{\text{test}}, \hat{\mathbf{Y}}_{\text{test}})$ 
  - leads to a beautiful theory of loss calibration

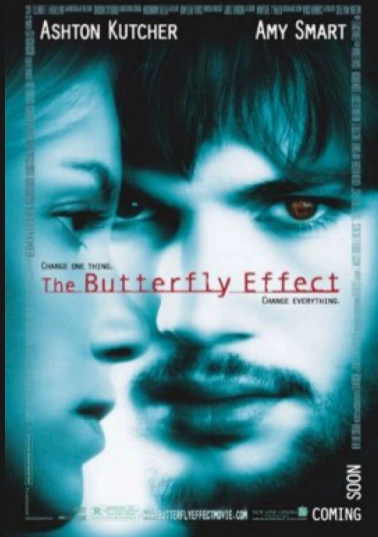
# Ill-conditioning

$$\mathbf{X} = \begin{bmatrix} 0 & \epsilon \\ 1 & 1 \end{bmatrix}, \quad \mathbf{y} = [1 \quad -1]$$

- Solving linear least squares regression:

$$\mathbf{w} = \mathbf{y}\mathbf{X}^{-1} = [1 \quad -1] \begin{bmatrix} -1/\epsilon & 1 \\ 1/\epsilon & 0 \end{bmatrix} = [-2/\epsilon \quad 1]$$

- Slight perturbation leads to chaotic behaviour!
- Happens whenever  $\mathbf{X}$  is ill-conditioned, i.e., (close to) rank deficient



# Tikhonov Regularization, a.k.a. Ridge Regression

$$\min_{\mathbf{W}} \frac{1}{n} \|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{W}\|_F^2$$

- Normal equation:  $\mathbf{W}(\mathbf{X}\mathbf{X}^\top + \lambda I) = \mathbf{Y}\mathbf{X}^\top$
- Regularization const.  $\lambda$  controls trade-off
  - $\lambda = 0$  reduces to ordinary linear regression
  - $\lambda = \infty$  reduces to  $\mathbf{W} \equiv \mathbf{0}$
  - intermediate  $\lambda$  restricts output to be  $\frac{1}{\lambda}$  proportional to input
- May choose to **not** regularize offset  $\mathbf{b}$



---

A. N. Tikhonov. "Solution of incorrectly formulated problems and the regularization method". *Soviet Mathematics*, vol. 4, no. 4 (1963), pp. 1035–1038, A. E. Hoerl and R. W. Kennard. "Ridge regression: Biased estimation for nonorthogonal problems". *Technometrics*, vol. 12, no. 1 (1970), pp. 55–67.

# Data Augmentation

$$\frac{1}{n} \|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_{\text{F}}^2 + \boxed{\lambda \|\mathbf{W}\|_{\text{F}}^2} = \frac{1}{n} \|\mathbf{W} \underbrace{\begin{bmatrix} \mathbf{X} & \sqrt{n\lambda}I \end{bmatrix}}_{\tilde{\mathbf{X}}} - \underbrace{\begin{bmatrix} \mathbf{Y} & \mathbf{0} \end{bmatrix}}_{\tilde{\mathbf{Y}}}\|_{\text{F}}^2$$

- Augment  $\mathbf{X}$  with  $\sqrt{n\lambda}I$ , i.e.  $p$  data points  $\mathbf{x}_j = \sqrt{n\lambda}\mathbf{e}_j, j = 1, \dots, p$
- Augment  $\mathbf{Y}$  with zero
- Shrinks  $\mathbf{W}$  towards origin

regularization = data augmentation

# Sparsity

- Regularization  $\iff$  constraint:

$$\min_{\|\mathbf{W}\|_F \leq \gamma} \frac{1}{n} \|\mathbf{WX} - \mathbf{Y}\|_F^2$$

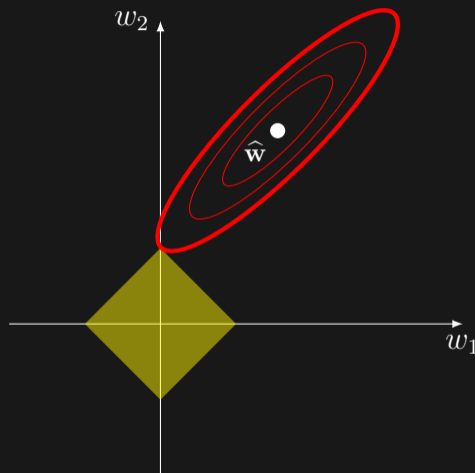
- Ridge regression  $\rightarrow$  dense  $\mathbf{W}$ 
  - more computation / communication
  - harder to interpret

- Lasso (Tibshirani, 1996):

$$\min_{\|\mathbf{W}\|_1 \leq \gamma} \frac{1}{n} \|\mathbf{WX} - \mathbf{Y}\|_F^2$$

- Regularization  $\iff$  constraint:

$$\min_{\mathbf{W}} \frac{1}{n} \|\mathbf{WX} - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{W}\|_1$$



# Task Regularization

$$\min_{\mathbf{W}} \frac{1}{n} \|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_{\text{F}}^2 + \lambda \|\mathbf{W}\|_{\text{F}}^2 \quad \equiv \quad \min_{\mathbf{w}_{\tau}} \frac{1}{n} \|\mathbf{w}_{\tau}\mathbf{X} - \mathbf{y}_{\tau}\|_{\text{F}}^2 + \lambda \|\mathbf{w}_{\tau}\|_2^2, \quad \forall \tau = 1, \dots, t$$

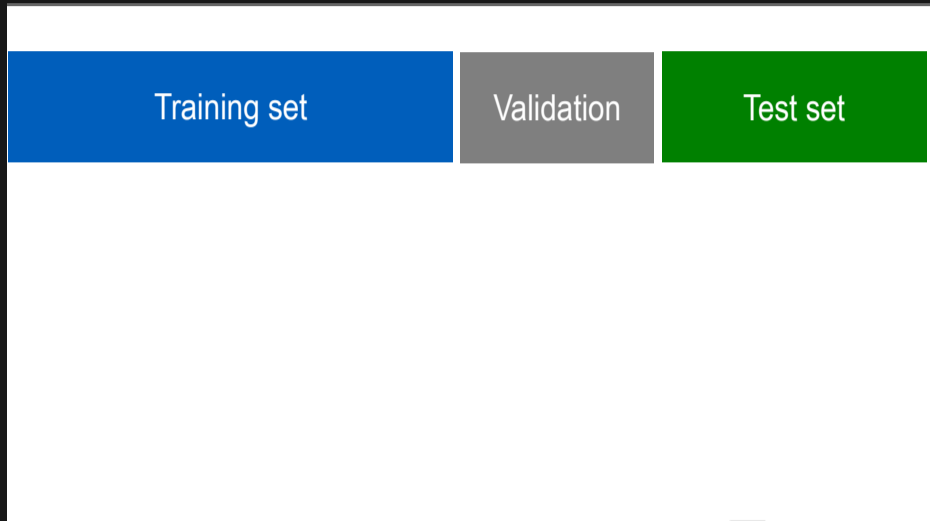
- In other words, the tasks are independent and can be solved separately
- Sometimes lumping tasks together (LHS) is computationally more efficient
- If tasks are related, can consider a kind of **low-rank regularization**:

$$\min_{\mathbf{W}} \frac{1}{n} \|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_{\text{F}}^2 + \lambda \|\mathbf{W}\|_{\text{tr}},$$

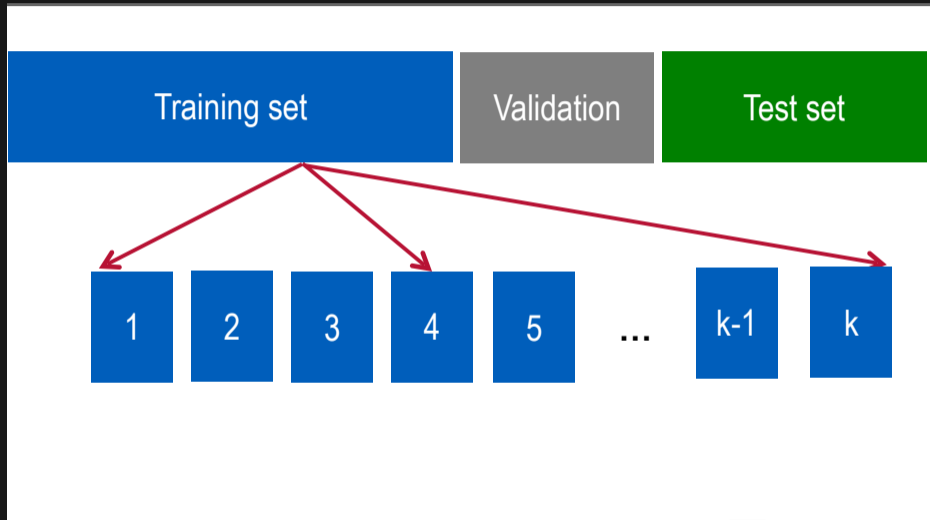
where  $\|A\|_{\text{tr}}$  is the sum of singular values (i.e., the **trace norm**).



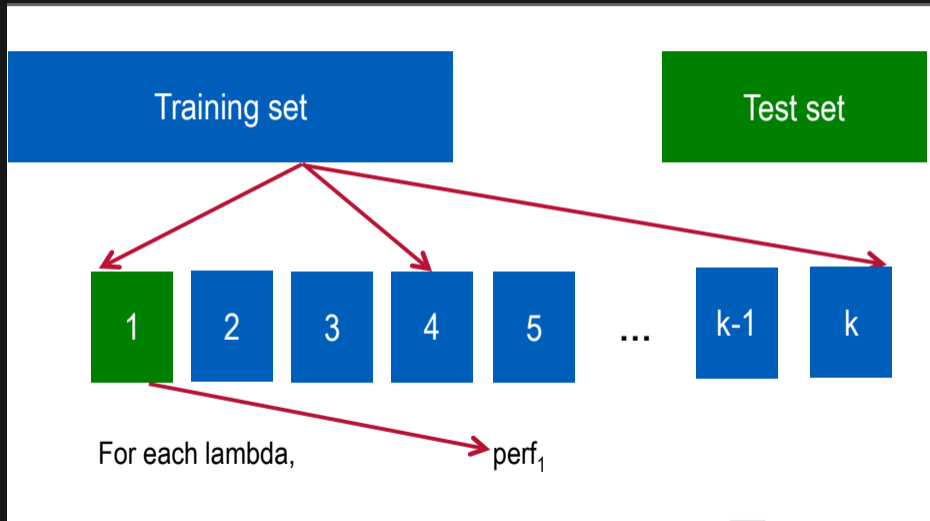
# Cross-validation



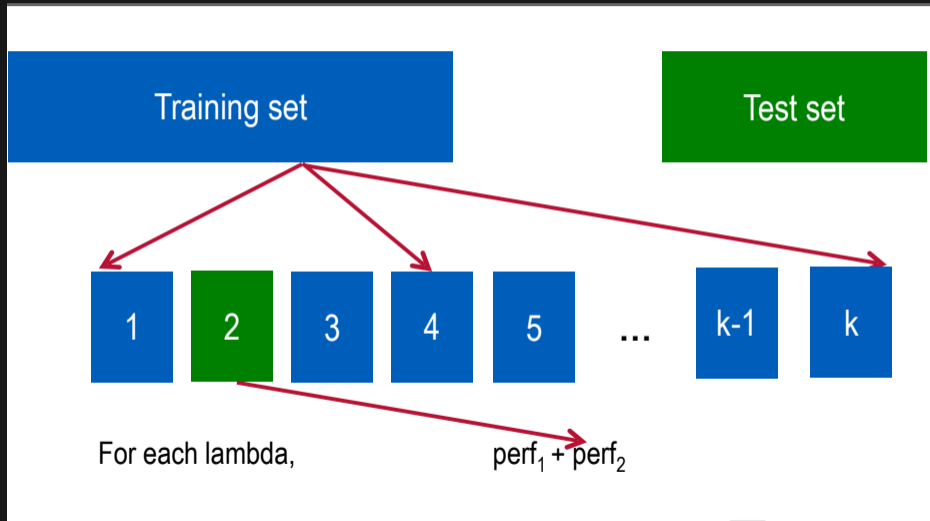
# Cross-validation



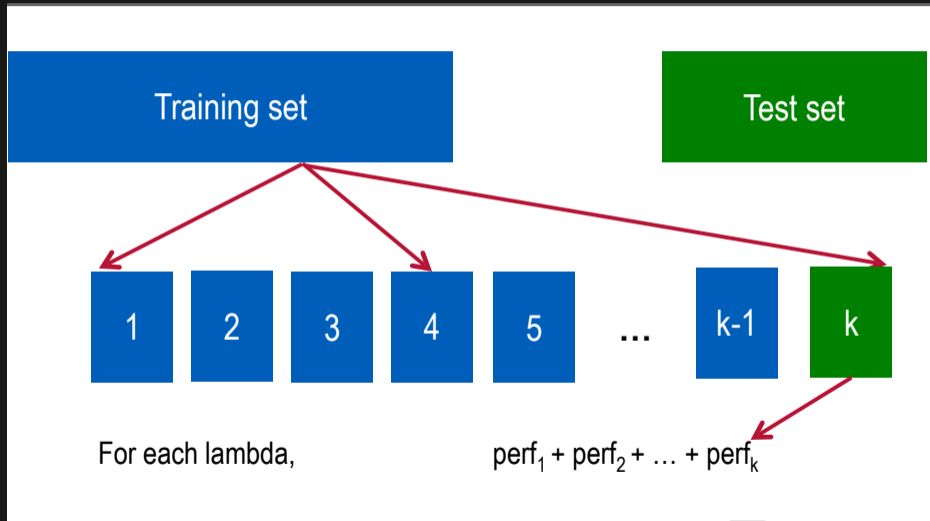
# Cross-validation



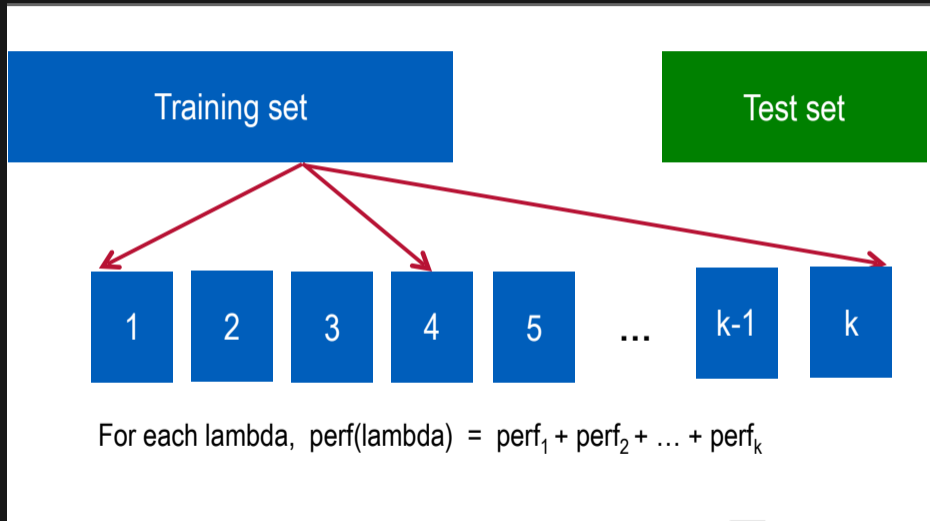
# Cross-validation



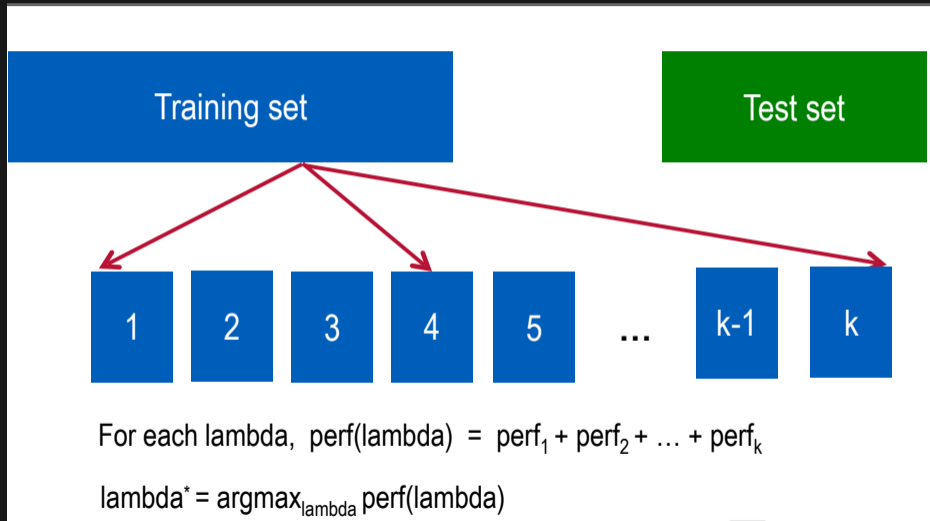
# Cross-validation



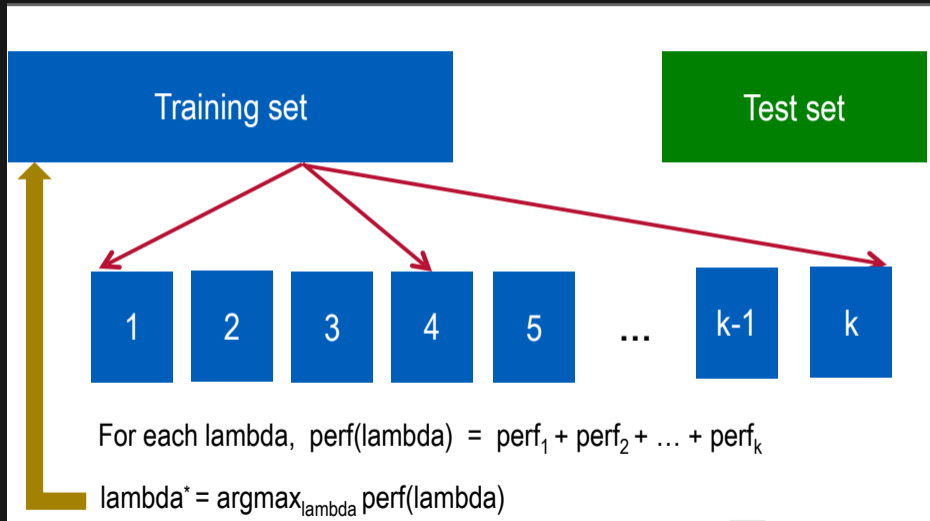
# Cross-validation



# Cross-validation



# Cross-validation





# Cross-validation

