# CS480/680: Introduction to Machine Learning
## Lec 02: Linear Regression

Yaoliang Yu
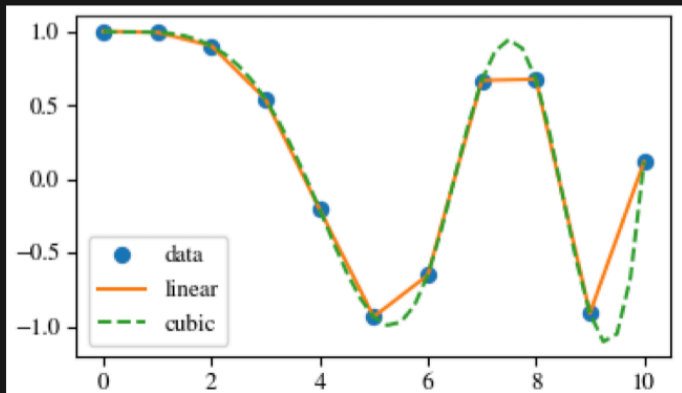
UNIVERSITY OF WATERLOO | FACULTY OF MATHEMATICS
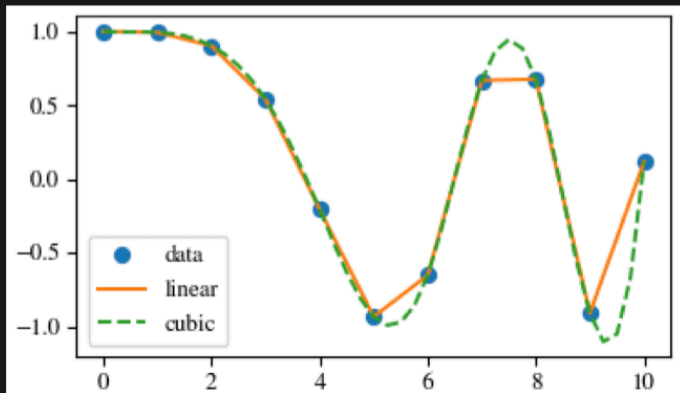DAVID R. CHERITON SCHOOL OF COMPUTER SCIENCE

May 13, 2024

# Regression

- Given training data $\{(\mathbf{x}_i, \mathbf{y}_i)\}$, find $f : \mathcal{X} \to \mathcal{Y}$ such that $f(\mathbf{x}_i) \approx y_i$

# Regression

- Given training data $\{(\mathbf{x}_i, \mathbf{y}_i)\}$, find $f : \mathcal{X} \to \mathcal{Y}$ 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 $t = \infty$

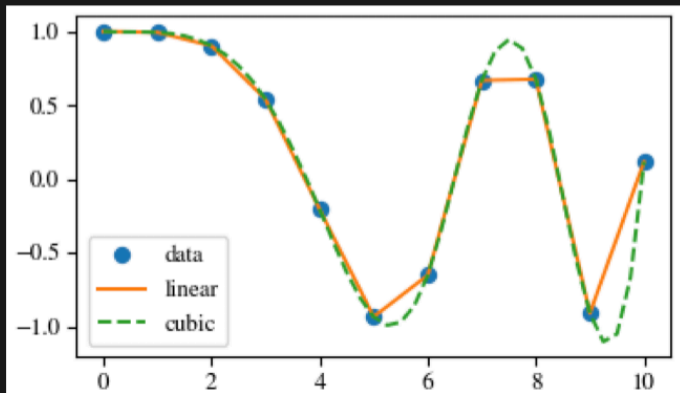# Regression

- Given training data $\{(\mathbf{x}_i, \mathbf{y}_i)\}$, find $f : \mathcal{X} \to \mathcal{Y}$ 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
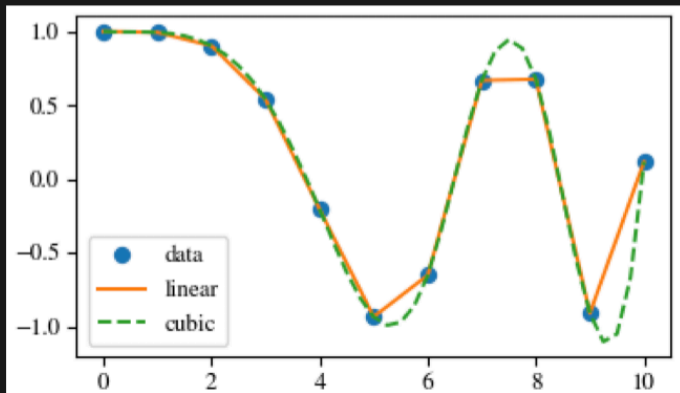  - $\mathbf{y}_i \in \mathcal{Y} \subseteq \mathbb{R}^t$: $t$ responses, e.g. $t = 1$ or even $t = \infty$

# Regression

- Given training data $\{(\mathbf{x}_i, \mathbf{y}_i)\}$, find $f : \mathcal{X} \to \mathcal{Y}$ 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



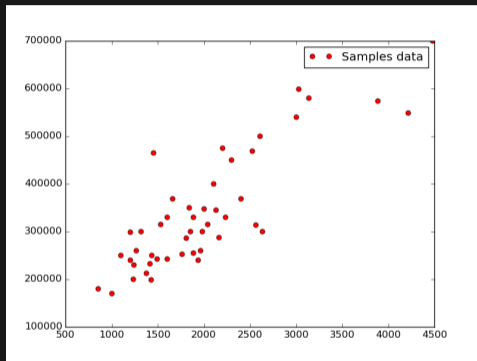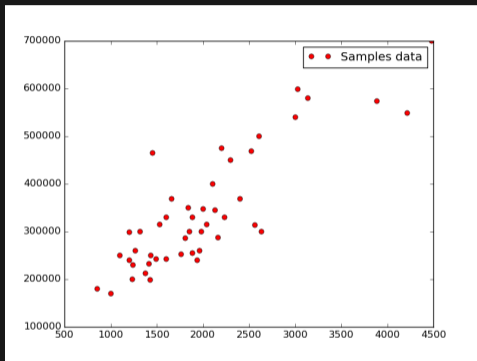- Prior knowledge on the functional form of $f$

- Linear vs. nonlinear

# Some Examples



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

# 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 $\{(\mathbf{x}_i, \mathbf{y}_i) : i = 1, \ldots, n\}$, there exist infinitely many functions $f$ such that for all $i$,
>
> $$f(\mathbf{x}_i) = \mathbf{y}_i.$$

- No amount of training data is enough to decide on a unique $f$!

- On new data $\mathbf{x}$, our prediction $\hat{\mathbf{y}} = f(\mathbf{x})$ can vary wildly!

- This is where prior knowledge of $f$ comes into play

- Occam's razor: "the simplest explanation is usually the correct one"

# The Difficulty

> **Theorem: Exact interpolation is always possible**
>
> For any[*] finite training data $\{(\mathbf{x}_i, \mathbf{y}_i) : i = 1, \ldots, n\}$, there exist infinitely many functions $f$ such that for all $i$,
>
> $$f(\mathbf{x}_i) = \mathbf{y}_i.$$

- No amount of training data is enough to decide on a unique $f$!

- On new data $\mathbf{x}$, our prediction $\hat{\mathbf{y}} = f(\mathbf{x})$ can vary wildly!

- This is where prior knowledge of $f$ comes into play

- Occam's razor: "the simplest explanation is usually the correct one"

# The Difficulty

> **Theorem: Exact interpolation is always possible**
>
> For any[*] finite training data $\{(\mathbf{x}_i, \mathbf{y}_i) : i = 1, \ldots, n\}$, there exist infinitely many functions $f$ such that for all $i$,
>
> $$f(\mathbf{x}_i) = \mathbf{y}_i.$$

- No amount of training data is enough to decide on a unique $f$!

- On new data $\mathbf{x}$, our prediction $\hat{\mathbf{y}} = f(\mathbf{x})$ can vary wildly!

- This is where prior knowledge of $f$ comes into play

- Occam's razor: "the simplest explanation is usually the correct one"

# The Difficulty

> **Theorem: Exact interpolation is always possible**
>
> For any[*] finite training data $\{(\mathbf{x}_i, \mathbf{y}_i) : i = 1, \ldots, n\}$, there exist infinitely many functions $f$ such that for all $i$,
>
> $$f(\mathbf{x}_i) = \mathbf{y}_i.$$

- No amount of training data is enough to decide on a unique $f$!

- On new data $\mathbf{x}$, our prediction $\hat{\mathbf{y}} = f(\mathbf{x})$ can vary wildly!

- This is where prior knowledge of $f$ comes into play

- Occam's razor: "the simplest explanation is usually the correct one"

# The Difficulty

**Theorem: Exact interpolation is always possible**

For any[*] finite training data $\{(\mathbf{x}_i, \mathbf{y}_i) : i = 1, \ldots, n\}$, there exist infinitely many functions $f$ such that for all $i$,

$$f(\mathbf{x}_i) = \mathbf{y}_i.$$

- No amount of training data is enough to decide on a unique $f$!

- On new data $\mathbf{x}$, our prediction $\hat{\mathbf{y}} = f(\mathbf{x})$ can vary wildly!

- This is where prior knowledge of $f$ comes into play

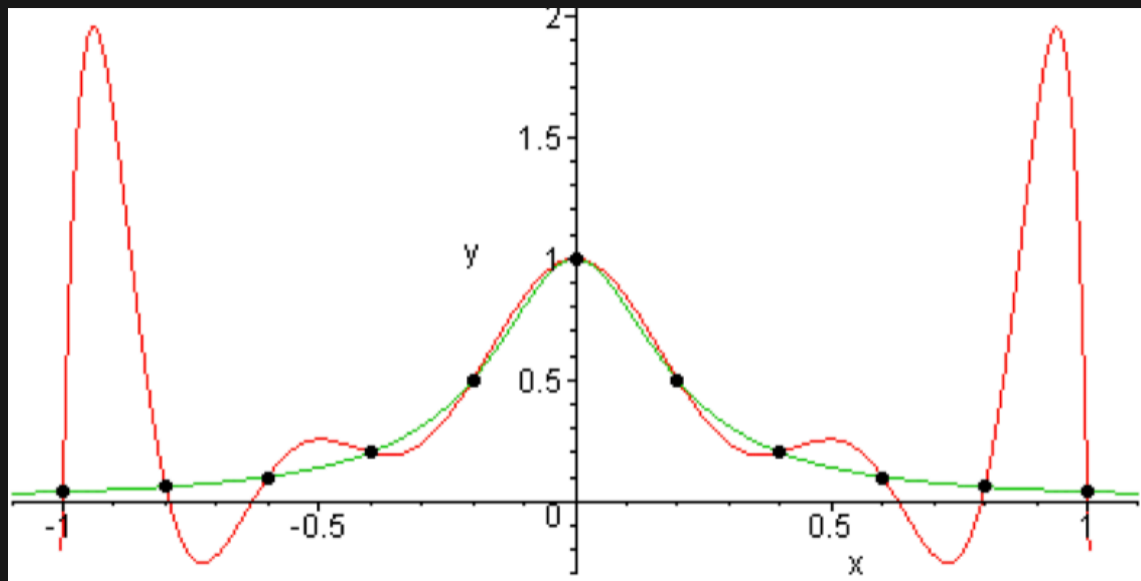- Occam's razor: "the simplest explanation is usually the correct one"

# Statistical Learning

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

- Least squares regression: $\min\limits_{f:\mathcal{X}\to\mathcal{Y}} \; \mathbb{E}\|f(\mathsf{X}) - \mathsf{Y}\|_2^2$

- Regression function: $m(\mathbf{x}) = \mathbb{E}[\mathsf{Y}|\mathsf{X} = \mathbf{x}]$

- Needs to know the distribution $\mathbb{P}$, i.e., all pairs $(\mathsf{X}, \mathsf{Y})$!

- Changing the square loss changes the regression function accordingly

# Statistical Learning

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

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

- Least squares regression: $\min\limits_{f:\mathcal{X} \to \mathcal{Y}} \ \mathbb{E}\|f(X) - Y\|_2^2$

- Regression function: $m(\mathbf{x}) = \mathbb{E}[Y|X = \mathbf{x}]$

- Needs to know the distribution $\mathbb{P}$, i.e., all pairs $(X, Y)$!

- Changing the square loss changes the regression function accordingly

# Statistical Learning

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

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

- Least squares regression: $\min\limits_{f: \mathcal{X} \to \mathcal{Y}} \; \mathbb{E}\|f(X) - Y\|_2^2$

- Regression function: $m(\mathbf{x}) = \mathbb{E}[Y|X = \mathbf{x}]$

- Needs to know the distribution $\mathbb{P}$, i.e., all pairs $(X, Y)$!

- Changing the square loss changes the regression function accordingly

# Statistical Learning

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

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

- Least squares regression: $\min_{f:\mathcal{X}\to\mathcal{Y}} \mathbb{E}\|f(X) - Y\|_2^2$

- Regression function: $m(\mathbf{x}) = \mathbb{E}[Y|X = \mathbf{x}]$

- Needs to know the distribution $\mathbb{P}$, i.e., all pairs $(X, Y)$!

- Changing the square loss changes the regression function accordingly

# Statistical Learning

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

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

- Least squares regression: $\min\limits_{f: \mathcal{X} \to \mathcal{Y}} \mathbb{E}\|f(X) - Y\|_2^2$

- Regression function: $m(\mathbf{x}) = \mathbb{E}[Y|X = \mathbf{x}]$

- Needs to know the distribution $\mathbb{P}$, i.e., all pairs $(X, Y)$!

- Changing the square loss changes the regression function accordingly

# Statistical Learning

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

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

- Least squares regression: $\min\limits_{f:\mathcal{X}\to\mathcal{Y}} \ \mathbb{E}\|f(X) - Y\|_2^2$

- Regression function: $m(\mathbf{x}) = \mathbb{E}[Y|X = \mathbf{x}]$

- Needs to know the distribution $\mathbb{P}$, i.e., all pairs $(X, Y)$!

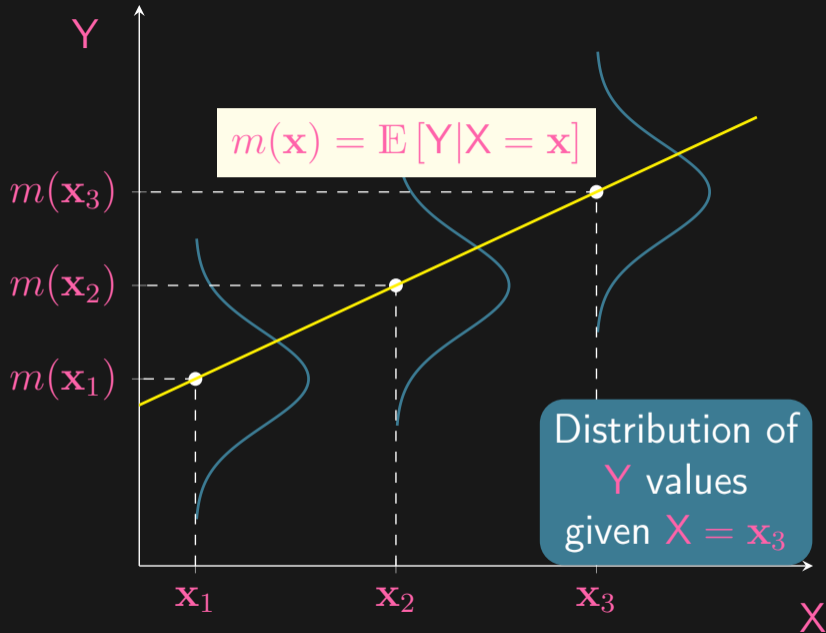- Changing the square loss changes the regression function accordingly

# Statistical Learning

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

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

- Least squares regression: $\min\limits_{f:\mathcal{X}\to\mathcal{Y}} \; \mathbb{E}\|f(X) - Y\|_2^2$

- Regression function: $m(\mathbf{x}) = \mathbb{E}[Y|X = \mathbf{x}]$

- Needs to know the distribution $\mathbb{P}$, i.e., all pairs $(X, Y)$!

- Changing the square loss changes the regression function accordingly

$m(\mathbf{x}) = \mathbb{E}\left[\mathsf{Y}|\mathsf{X}=\mathbf{x}\right]$

$m(\mathbf{x}_3)$

$m(\mathbf{x}_2)$

$m(\mathbf{x}_1)$

Distribution of $\mathsf{Y}$ values given $\mathsf{X} = \mathbf{x}_3$

$\mathbf{x}_1$ $\quad$ $\mathbf{x}_2$ $\quad$ $\mathbf{x}_3$

X

Y

# Bias-Variance Decomposition

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

- 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

# Bias-Variance Decomposition

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

- 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

# Bias-Variance Decomposition

$$\mathbb{E}\|f(\mathsf{X}) - \mathsf{Y}\|_2^2 = \mathbb{E}\|f(\mathsf{X}) - m(\mathsf{X}) + m(\mathsf{X}) - \mathsf{Y}\|_2^2$$

$$= \mathbb{E}\|f(\mathsf{X}) - m(\mathsf{X})\|_2^2 + \mathbb{E}\|m(\mathsf{X}) - \mathsf{Y}\|_2^2$$

$$\underbrace{+2\mathbb{E}\langle f(\mathsf{X}) - m(\mathsf{X}), m(\mathsf{X}) - \mathsf{Y}\rangle}$$

$$= \underbrace{\mathbb{E}\|f(\mathsf{X}) - m(\mathsf{X})\|_2^2}_{\text{bias}^2} + \underbrace{\mathbb{E}\|m(\mathsf{X}) - \mathsf{Y}\|_2^2}_{\text{noise variance}}$$

- 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

# Bias-Variance Decomposition

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

- 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 → Training

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

- Replace expectation with sample average: $(\mathsf{X}_i, \mathsf{Y}_i) \overset{i.i.d.}{\sim} P$

- Finite training set → exact interpolation paradox!

- Need to restrict the form of $f$, using prior knowledge

- (Uniform) law of large numbers: as training data size $n \to \infty$,

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

# Sampling → Training

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

- Replace expectation with sample average: $(\mathsf{X}_i, \mathsf{Y}_i) \overset{i.i.d.}{\sim} P$

- Finite training set → exact interpolation paradox!

- Need to restrict the form of $f$, using prior knowledge

- (Uniform) law of large numbers: as training data size $n \to \infty$,

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

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

- Replace expectation with sample average: $(\mathsf{X}_i, \mathsf{Y}_i) \overset{i.i.d.}{\sim} P$

- Finite training set → exact interpolation paradox!

- Need to restrict the form of $f$, using prior knowledge

- (Uniform) law of large numbers: as training data size $n \to \infty$,

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

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

- Replace expectation with sample average: $(\mathsf{X}_i, \mathsf{Y}_i) \overset{i.i.d.}{\sim} P$

- Finite training set $\to$ exact interpolation paradox!

- Need to restrict the form of $f$, using prior knowledge

- (Uniform) law of large numbers: as training data size $n \to \infty$,

  $\hat{\mathbb{E}} \to \mathbb{E}$ and (hopefully) $\arg\min \hat{\mathbb{E}} \to \arg\min \mathbb{E}$

# Sampling → Training

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

- Replace expectation with sample average: $(\mathsf{X}_i, \mathsf{Y}_i) \overset{i.i.d.}{\sim} P$

- Finite training set → exact interpolation paradox!

- Need to restrict the form of $f$, using prior knowledge

- (Uniform) law of large numbers: as training data size $n \to \infty$,

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

# 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: $\mathbf{x} \leftarrow \binom{\mathbf{x}}{1}$, $\mathbf{W} \leftarrow [W, \mathbf{b}]$, hence $f(\mathbf{x}) = \mathbf{W}\mathbf{x}$

- In matrix form: $\frac{1}{n} \sum_i \|f(\mathbf{x}_i) - \mathbf{y}_i\|_2^2 \;=\; \frac{1}{n} \|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_F^2$

S. M. Stigler. "Gauss and the Invention of Least Squares". *The Annals of Statistics*, vol. 9, no. 3 (1981), pp. 465–474.

# 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: $\mathbf{x} \leftarrow \binom{\mathbf{x}}{1}$, $\mathbf{W} \leftarrow [W, \mathbf{b}]$, hence $f(\mathbf{x}) = \mathbf{W}\mathbf{x}$

- In matrix form: $\frac{1}{n} \sum_i \|f(\mathbf{x}_i) - \mathbf{y}_i\|_2^2 \;=\; \frac{1}{n} \|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_F^2$

S. M. Stigler. "Gauss and the Invention of Least Squares". *The Annals of Statistics*, vol. 9, no. 3 (1981), pp. 465–474.

# 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: $\mathbf{x} \leftarrow \binom{\mathbf{x}}{1}$, $\mathbf{W} \leftarrow [W, \mathbf{b}]$, hence $f(\mathbf{x}) = \mathbf{W}\mathbf{x}$

- In matrix form: $\frac{1}{n}\sum_i \|f(\mathbf{x}_i) - \mathbf{y}_i\|_2^2 \;=\; \frac{1}{n}\|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_F^2$

S. M. Stigler. "Gauss and the Invention of Least Squares". *The Annals of Statistics*, vol. 9, no. 3 (1981), pp. 465–474.

# 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: $\frac{1}{n} \sum_i \|f(\mathbf{x}_i) - \mathbf{y}_i\|_2^2 \quad = \quad \frac{1}{n} \|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_F^2$

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

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

S. M. Stigler. "Gauss and the Invention of Least Squares". *The Annals of Statistics*, vol. 9, no. 3 (1981), pp. 465–474.

# 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: $\mathbf{x} \leftarrow \binom{\mathbf{x}}{1}$, $\mathbf{W} \leftarrow [W, \mathbf{b}]$, hence $f(\mathbf{x}) = \mathbf{W}\mathbf{x}$

- In matrix form: $\frac{1}{n} \sum_i \|f(\mathbf{x}_i) - \mathbf{y}_i\|_2^2 \;\; = \;\; \frac{1}{n} \|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_F^2$

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

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

---

S. M. Stigler. "Gauss and the Invention of Least Squares". *The Annals of Statistics*, vol. 9, no. 3 (1981), pp. 465–474.

# 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: $\mathbf{x} \leftarrow \binom{\mathbf{x}}{1}$, $\mathbf{W} \leftarrow [W, \mathbf{b}]$, hence $f(\mathbf{x}) = \mathbf{W}\mathbf{x}$

- In matrix form: $\frac{1}{n} \sum_i \|f(\mathbf{x}_i) - \mathbf{y}_i\|_2^2 \quad = \quad \frac{1}{n}\|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_\mathsf{F}^2$
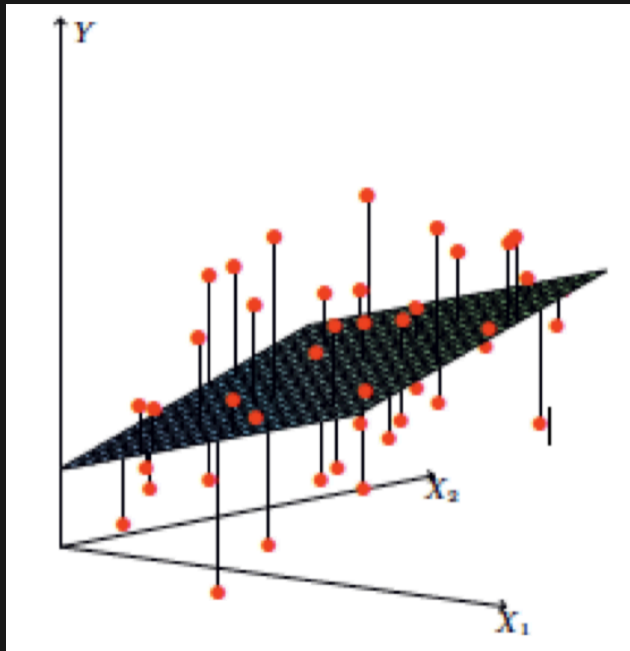
  - $\mathbf{X} = [\mathbf{x}_1, \ldots, \mathbf{x}_n] \in \mathbb{R}^{(d+1) \times n}$, $\mathbf{Y} = [\mathbf{y}_1, \ldots, \mathbf{y}_n] \in \mathbb{R}^{t \times n}$

  - $\|A\|_\mathsf{F} = \sqrt{\sum_{ij} a_{ij}^2}$

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

S. M. Stigler. "Gauss and the Invention of Least Squares". *The Annals of Statistics*, vol. 9, no. 3 (1981), pp. 465–474.

# Calculus Detour

- Let $f : \mathbb{R}^p \to \mathbb{R}$ be a smooth real-valued function

- Fix an inner product $\langle \cdot, \cdot \rangle$

- Define the gradient $\nabla f : \mathbb{R}^p \to \mathbb{R}^p$ as

$$\frac{\mathrm{d}f(\mathbf{w} + t\mathbf{z})}{\mathrm{d}t} \bigg|_{t=0} = \langle \nabla f(\mathbf{w}), \mathbf{z} \rangle$$

  - LHS = the usual 1-variable derivative of the composite function $t \mapsto f(\mathbf{w} + t\mathbf{z})$

  - $\mathbf{w}$ and $\mathbf{z}$ are fixed as constants

  - gradient $\nabla f$ is equivalent map of directional derivative over the inner product gin space

- Chain rule still holds

# Calculus Detour

- Let $f : \mathbb{R}^p \to \mathbb{R}$ be a smooth real-valued function

- Fix an inner product $\langle \cdot, \cdot \rangle$

- Define the gradient $\nabla f : \mathbb{R}^p \to \mathbb{R}^p$ as

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

  - LHS = the usual 1-variable derivative of the univariate function $t \mapsto f(\mathbf{w} + t\mathbf{z})$ with $\mathbf{w}$ and $\mathbf{z}$ are fixed as constants

  - gradient $\nabla f$ is represent map of directional derivative over the inner product-set space

- Chain rule still holds

# Calculus Detour

- Let $f : \mathbb{R}^p \to \mathbb{R}$ be a smooth real-valued function

- Fix an inner product $\langle \cdot, \cdot \rangle$

- Define the gradient $\nabla f : \mathbb{R}^p \to \mathbb{R}^p$ as

$$\frac{\mathrm{d}f(\mathbf{w} + t\mathbf{z})}{\mathrm{d}t} \big\rvert_{t=0} = \langle \nabla f(\mathbf{w}), \mathbf{z} \rangle$$

- LHS = the usual 1-dim'l derivative of the univariate function $t \mapsto f(\mathbf{w} + t\mathbf{z})$ with $\mathbf{w}$ and $\mathbf{z}$ are fixed as constants

- $\Rightarrow$ gradient $\nabla f$ is representation of directional derivative over the inner product we choose

- Chain rule still holds

# Calculus Detour

- Let $f : \mathbb{R}^p \to \mathbb{R}$ be a smooth real-valued function

- Fix an inner product $\langle \cdot, \cdot \rangle$

- Define the gradient $\nabla f : \mathbb{R}^p \to \mathbb{R}^p$ as

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

  - 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 the inner product we choose

- Chain rule still holds

# Calculus Detour

- Let $f : \mathbb{R}^p \to \mathbb{R}$ be a smooth real-valued function

- Fix an inner product $\langle \cdot, \cdot \rangle$

- Define the gradient $\nabla f : \mathbb{R}^p \to \mathbb{R}^p$ as

$$\frac{\mathrm{d} f(\mathbf{w} + t\mathbf{z})}{\mathrm{d}t} \Big\rvert_{t=0} = \langle \nabla f(\mathbf{w}), \mathbf{z} \rangle$$

  - 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 the inner product we choose

- Chain rule still holds

# Calculus Detour

- Let $f : \mathbb{R}^p \to \mathbb{R}$ be a smooth real-valued function

- Fix an inner product $\langle \cdot, \cdot \rangle$

- Define the gradient $\nabla f : \mathbb{R}^p \to \mathbb{R}^p$ as

$$\frac{\mathrm{d} f(\mathbf{w} + t\mathbf{z})}{\mathrm{d}t} \big|_{t=0} = \langle \nabla f(\mathbf{w}), \mathbf{z} \rangle$$

  - 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 the inner product we choose

- Chain rule still holds

# Calculus Detour

- Let $f : \mathbb{R}^p \to \mathbb{R}$ be a smooth real-valued function

- Fix an inner product $\langle \cdot, \cdot \rangle$

- Define the gradient $\nabla f : \mathbb{R}^p \to \mathbb{R}^p$ as

$$\frac{\mathrm{d}f(\mathbf{w} + t\mathbf{z})}{\mathrm{d}t} \big\vert_{t=0} = \langle \nabla f(\mathbf{w}), \mathbf{z} \rangle$$

  - 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 the inner product we choose

- Chain rule still holds

# Calculus Detour

- Let $f : \mathbb{R}^p \to \mathbb{R}$ be a smooth real-valued function

- Fix an inner product $\langle \cdot, \cdot \rangle$

- Define the gradient $\nabla f : \mathbb{R}^p \to \mathbb{R}^p$ as

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

  – 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 the inner product we choose

- Chain rule still holds

## Example: Univariate functions

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

$$\frac{\mathrm{d}f(w + tz)}{\mathrm{d}t} \upharpoonright_{t=0} = f'(w + tz)z \upharpoonright_{t=0} = f'(w)z = \langle f'(w), z \rangle,$$

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

## Example: Partial derivatives

Consider $f : \mathbb{R}^p \to \mathbb{R}$ 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):

$$\frac{\mathrm{d}f(\mathbf{w} + t\mathbf{e}_j)}{\mathrm{d}t} \upharpoonright_{t=0} = \partial_j f(\mathbf{w}) = \langle \nabla f(\mathbf{w}), \mathbf{e}_j \rangle = [\nabla f(\mathbf{w})]_j,$$

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

### Example: Quadratic function

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

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

$$\frac{\mathrm{d}f(\mathbf{w} + t\mathbf{z})}{\mathrm{d}t} \upharpoonright_{t=0} = \langle \mathbf{w}, A\mathbf{z} \rangle + \langle \mathbf{z}, A\mathbf{w} + \mathbf{b} \rangle = \left\langle A^\top \mathbf{w} + A\mathbf{w} + \mathbf{b}, \mathbf{z} \right\rangle,$$

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

- $\langle \mathbf{a} + \mathbf{b}, \mathbf{x} + \mathbf{y} \rangle = \langle \mathbf{a}, \mathbf{x} \rangle + \langle \mathbf{a}, \mathbf{y} \rangle + \langle \mathbf{b}, \mathbf{x} \rangle + \langle \mathbf{b}, \mathbf{y} \rangle$

- $\langle \mathbf{a}, t\mathbf{b} \rangle = \langle t\mathbf{a}, \mathbf{b} \rangle = t \langle \mathbf{a}, \mathbf{b} \rangle$

- $\langle \mathbf{w}, A\mathbf{z} \rangle = \left\langle A^\top \mathbf{w}, \mathbf{z} \right\rangle, \ \langle A\mathbf{w}, \mathbf{z} \rangle = \left\langle \mathbf{w}, A^\top \mathbf{z} \right\rangle$

### Example: Quadratic function

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

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

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

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

- $\langle \mathbf{a} + \mathbf{b}, \mathbf{x} + \mathbf{y} \rangle = \langle \mathbf{a}, \mathbf{x} \rangle + \langle \mathbf{a}, \mathbf{y} \rangle + \langle \mathbf{b}, \mathbf{x} \rangle + \langle \mathbf{b}, \mathbf{y} \rangle$

- $\langle \mathbf{a}, t\mathbf{b} \rangle = \langle t\mathbf{a}, \mathbf{b} \rangle = t \langle \mathbf{a}, \mathbf{b} \rangle$

- $\langle \mathbf{w}, A\mathbf{z} \rangle = \langle A^\top \mathbf{w}, \mathbf{z} \rangle, \ \langle A\mathbf{w}, \mathbf{z} \rangle = \langle \mathbf{w}, A^\top \mathbf{z} \rangle$

### Example: Quadratic function

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

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

$$\frac{\mathrm{d} f(\mathbf{w} + t\mathbf{z})}{\mathrm{d} t} \restriction_{t=0} = \langle \mathbf{w}, A\mathbf{z} \rangle + \langle \mathbf{z}, A\mathbf{w} + \mathbf{b} \rangle = \langle A^\top \mathbf{w} + A\mathbf{w} + \mathbf{b}, \mathbf{z} \rangle,$$

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

- $\langle \mathbf{a} + \mathbf{b}, \mathbf{x} + \mathbf{y} \rangle = \langle \mathbf{a}, \mathbf{x} \rangle + \langle \mathbf{a}, \mathbf{y} \rangle + \langle \mathbf{b}, \mathbf{x} \rangle + \langle \mathbf{b}, \mathbf{y} \rangle$

- $\langle \mathbf{a}, t\mathbf{b} \rangle = \langle t\mathbf{a}, \mathbf{b} \rangle = t \langle \mathbf{a}, \mathbf{b} \rangle$

- $\langle \mathbf{w}, A\mathbf{z} \rangle = \langle A^\top \mathbf{w}, \mathbf{z} \rangle, \ \langle A\mathbf{w}, \mathbf{z} \rangle = \langle \mathbf{w}, A^\top \mathbf{z} \rangle$
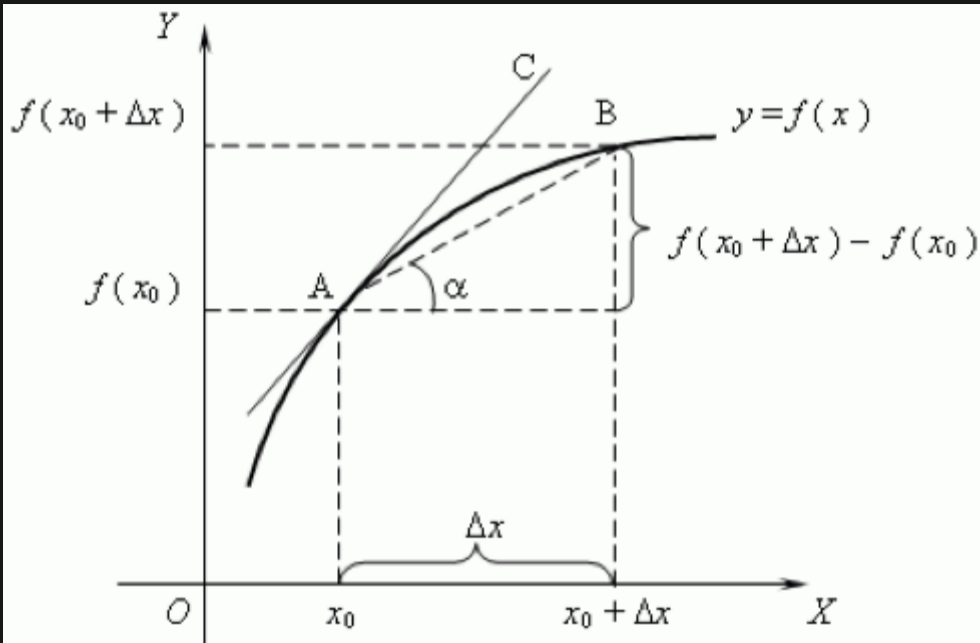
### Example: Quadratic function

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

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

$$\frac{\mathrm{d}f(\mathbf{w} + t\mathbf{z})}{\mathrm{d}t} \Big\vert_{t=0} = \langle \mathbf{w}, A\mathbf{z} \rangle + \langle \mathbf{z}, A\mathbf{w} + \mathbf{b} \rangle = \left\langle A^\top \mathbf{w} + A\mathbf{w} + \mathbf{b}, \mathbf{z} \right\rangle,$$

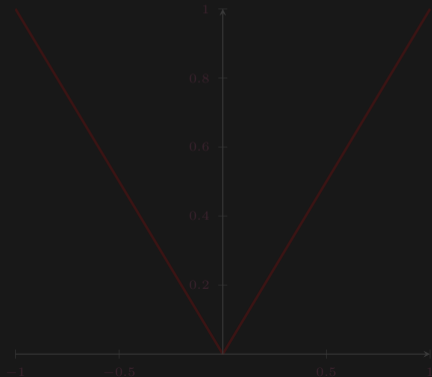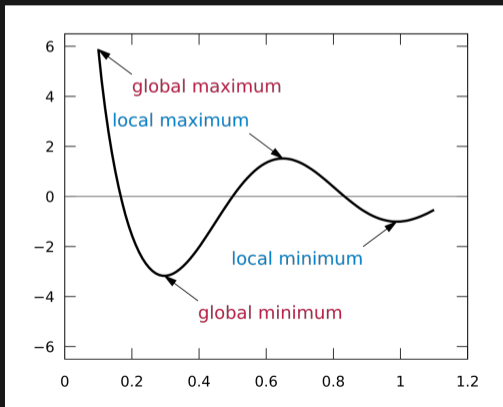i.e., $\boxed{\nabla f(\mathbf{w}) = (A^\top + A)\mathbf{w} + \mathbf{b}}$.

- $\langle \mathbf{a} + \mathbf{b}, \mathbf{x} + \mathbf{y} \rangle = \langle \mathbf{a}, \mathbf{x} \rangle + \langle \mathbf{a}, \mathbf{y} \rangle + \langle \mathbf{b}, \mathbf{x} \rangle + \langle \mathbf{b}, \mathbf{y} \rangle$

- $\langle \mathbf{a}, t\mathbf{b} \rangle = \langle t\mathbf{a}, \mathbf{b} \rangle = t \langle \mathbf{a}, \mathbf{b} \rangle$

- $\langle \mathbf{w}, A\mathbf{z} \rangle = \left\langle A^\top \mathbf{w}, \mathbf{z} \right\rangle$, $\langle A\mathbf{w}, \mathbf{z} \rangle = \left\langle \mathbf{w}, A^\top \mathbf{z} \right\rangle$

# Optimality Condition

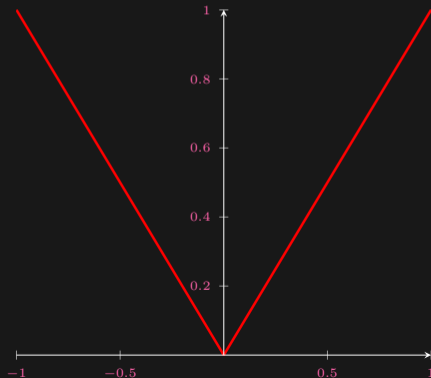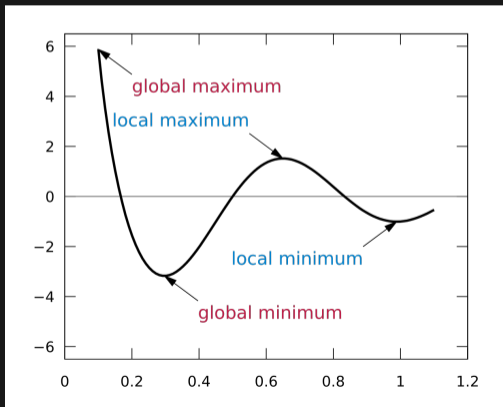**Theorem: Fermat's necessary condition for extremity**

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

# 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

$$\|\mathbf{WX} - \mathbf{Y}\|_{\mathsf{F}}^2 = \langle \mathbf{WX} - \mathbf{Y}, \mathbf{WX} - \mathbf{Y} \rangle$$
$$= \langle \mathbf{W}, \mathbf{WXX}^\top - 2\mathbf{YX}^\top \rangle + \langle \mathbf{Y}, \mathbf{Y} \rangle$$

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

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

- $\mathbf{X} \in \mathbb{R}^{(d+1) \times n}$ 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

# Solving Linear Regression

$$\|\mathbf{WX} - \mathbf{Y}\|_{\mathsf{F}}^2 = \langle \mathbf{WX} - \mathbf{Y}, \mathbf{WX} - \mathbf{Y} \rangle$$
$$= \langle \mathbf{W}, \mathbf{WXX}^\top - 2\mathbf{YX}^\top \rangle + \langle \mathbf{Y}, \mathbf{Y} \rangle$$

- 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

# Solving Linear Regression

$$\|\mathbf{WX} - \mathbf{Y}\|_{\mathsf{F}}^2 = \langle \mathbf{WX} - \mathbf{Y}, \mathbf{WX} - \mathbf{Y} \rangle$$
$$= \langle \mathbf{W}, \mathbf{WXX}^\top - 2\mathbf{YX}^\top \rangle + \langle \mathbf{Y}, \mathbf{Y} \rangle$$

- 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 $n \le d+1$, but a solution always exists

- Even when invertible, never compute the inverse directly!

- Instead, solve the linear system or apply iterative gradient algorithm

# Solving Linear Regression

$$\|\mathbf{WX} - \mathbf{Y}\|_\mathsf{F}^2 = \langle \mathbf{WX} - \mathbf{Y}, \mathbf{WX} - \mathbf{Y} \rangle$$
$$= \langle \mathbf{W}, \mathbf{WXX}^\top - 2\mathbf{YX}^\top \rangle + \langle \mathbf{Y}, \mathbf{Y} \rangle$$

- 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

# Solving Linear Regression

$$\|\mathbf{WX} - \mathbf{Y}\|_{\mathsf{F}}^2 = \langle \mathbf{WX} - \mathbf{Y}, \mathbf{WX} - \mathbf{Y} \rangle$$
$$= \langle \mathbf{W}, \mathbf{WXX}^\top - 2\mathbf{YX}^\top \rangle + \langle \mathbf{Y}, \mathbf{Y} \rangle$$

- 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:

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

- We may compare to the training error:

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

- Minimizing the training error as a means to reduce the test error

- Sometimes we even evaluate the test error using a different loss $\mathbb{L}(\mathbf{Y}_{\text{test}}, \hat{\mathbf{Y}}_{\text{test}})$

# 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:

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

- We may compare to the training error:

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

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

- 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:

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

- We may compare to the training error:

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

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

# 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:

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

- We may compare to the training error:

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

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

# 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:

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

- We may compare to the training error:

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

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

# 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:

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

- We may compare to the training error:

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

- 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

# 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:

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

- We may compare to the training error:

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

- 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

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

- Solving linear least squares regression:

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

- Slight perturbation leads to chaotic behaviour!

- Happens whenever $\mathbf{X}$ is ill-conditioned, i.e., (close to) rank deficient

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

- Solving linear least squares regression:

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

- Slight perturbation leads to chaotic behaviour!

- Happens whenever $\mathbf{X}$ is ill-conditioned, i.e., (close to) rank deficient

# Ill-conditioning



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

- Solving linear least squares regression:

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

- Slight perturbation leads to chaotic behaviour!

- Happens whenever $\mathbf{X}$ is ill-conditioned, i.e., (close to) rank deficient

# Ill-conditioning
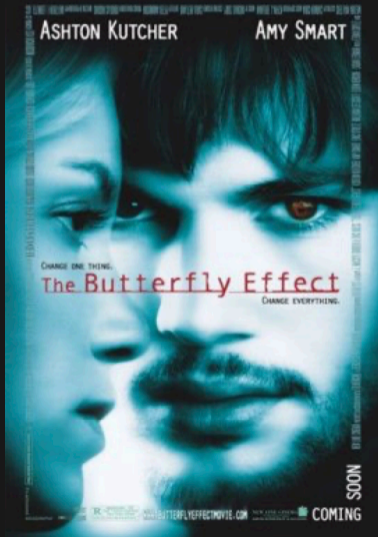
$$\mathbf{X} = \begin{bmatrix} 0 & \epsilon \\ 1 & 1 \end{bmatrix}, \qquad \mathbf{y} = \begin{bmatrix} 1 & -1 \end{bmatrix}$$

- Solving linear least squares regression:

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

- Slight perturbation leads to chaotic behaviour!

- Happens whenever **X** is ill-conditioned, i.e., (close to) rank deficient

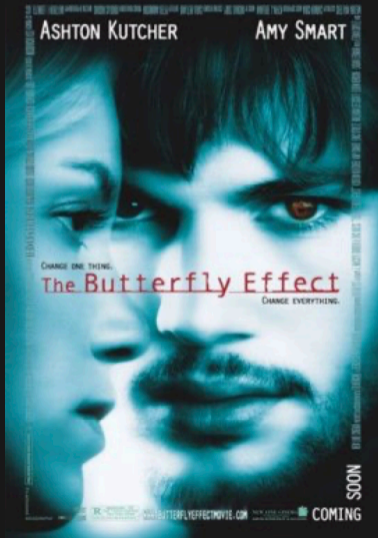# Tikhonov Regularization, a.k.a. Ridge Regression

$$\min_{\mathbf{W}} \quad \frac{1}{n}\|\mathbf{WX} - \mathbf{Y}\|_{\mathsf{F}}^2 + \boxed{\lambda\|\mathbf{W}\|_{\mathsf{F}}^2}$$

- Normal equation: $\mathbf{W}(\mathbf{XX}^\top + \lambda I) = \mathbf{YX}^\top$
- Regularization const. $\lambda$ controls trade-off

- 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.

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

$$\min_{\mathbf{W}} \quad \tfrac{1}{n}\|\mathbf{WX} - \mathbf{Y}\|_{\mathsf{F}}^2 + \boxed{\lambda\|\mathbf{W}\|_{\mathsf{F}}^2}$$

- Normal equation: $\mathbf{W}(\mathbf{XX}^\top + \lambda I) = \mathbf{YX}^\top$

- Regularization const. $\lambda$ controls trade-off

- 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.

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

$$\min_{\mathbf{W}} \ \frac{1}{n}\|\mathbf{WX} - \mathbf{Y}\|_{\mathsf{F}}^2 + \boxed{\lambda\|\mathbf{W}\|_{\mathsf{F}}^2}$$



- Normal equation: $\mathbf{W}(\mathbf{XX}^\top + \lambda I) = \mathbf{YX}^\top$
- 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 $b$

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.

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

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



- 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 $\mathbf{W} \equiv 0$
  - 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.

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

$$\min_{\mathbf{W}} \quad \frac{1}{n}\|\mathbf{WX} - \mathbf{Y}\|_{\mathsf{F}}^2 + \boxed{\lambda\|\mathbf{W}\|_{\mathsf{F}}^2}$$



- Normal equation: $\mathbf{W}(\mathbf{XX}^\top + \lambda I) = \mathbf{YX}^\top$
- Regularization const. $\lambda$ controls trade-off
    - $\lambda = 0$ reduces to ordinary linear regression
    - $\lambda = \infty$ reduces to $\mathbf{W} \equiv 0$
    - 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.

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

$$\min_{\mathbf{W}} \quad \frac{1}{n}\|\mathbf{WX} - \mathbf{Y}\|_{\mathsf{F}}^2 + \boxed{\lambda\|\mathbf{W}\|_{\mathsf{F}}^2}$$

- Normal equation: $\mathbf{W}(\mathbf{XX}^\top + \lambda I) = \mathbf{YX}^\top$
- 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.

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

$$\min_{\mathbf{W}} \ \frac{1}{n}\|\mathbf{WX} - \mathbf{Y}\|_{\mathsf{F}}^2 + \boxed{\lambda\|\mathbf{W}\|_{\mathsf{F}}^2}$$



- Normal equation: $\mathbf{W}(\mathbf{XX}^\top + \lambda I) = \mathbf{YX}^\top$
- 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

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

- 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

# Data Augmentation

$$\frac{1}{n}\|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_{\mathsf{F}}^2 + \boxed{\lambda\|\mathbf{W}\|_{\mathsf{F}}^2} = \frac{1}{n}\|\mathbf{W}\underbrace{\left[\mathbf{X} \quad \sqrt{n\lambda}I\right]}_{\tilde{\mathbf{X}}} - \underbrace{\left[\mathbf{Y} \quad \mathbf{0}\right]}_{\tilde{\mathbf{Y}}}\|_{\mathsf{F}}^2$$

- Augment $\mathbf{X}$ with $\sqrt{n\lambda}I$, i.e. $p$ data points $\mathbf{x}_j = \sqrt{n\lambda}\mathbf{e}_j, j = 1, \ldots, p$

- Augment $\mathbf{Y}$ with zero

- Shrinks $\mathbf{W}$ towards origin

# Data Augmentation

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

- Augment $\mathbf{X}$ with $\sqrt{n\lambda}I$, i.e. $p$ data points $\mathbf{x}_j = \sqrt{n\lambda}\mathbf{e}_j, j = 1, \ldots, p$

- Augment $\mathbf{Y}$ with zero

- Shrinks $\mathbf{W}$ towards origin

## Data Augmentation

$$\frac{1}{n}\|\mathbf{W}\mathbf{X} - \mathbf{Y}\|_{\mathsf{F}}^2 + \boxed{\lambda\|\mathbf{W}\|_{\mathsf{F}}^2} = \frac{1}{n}\|\mathbf{W}\underbrace{[\mathbf{X} \quad \sqrt{n\lambda}I]}_{\tilde{\mathbf{X}}} - \underbrace{[\mathbf{Y} \quad \mathbf{0}]}_{\tilde{\mathbf{Y}}}\|_{\mathsf{F}}^2$$

- Augment $\mathbf{X}$ with $\sqrt{n\lambda}I$, i.e. $p$ data points $\mathbf{x}_j = \sqrt{n\lambda}\mathbf{e}_j, j = 1, \ldots, p$

- Augment $\mathbf{Y}$ with zero

- Shrinks $\mathbf{W}$ towards origin

$$\boxed{\text{regularization} = \text{data augmentation}}$$

# Sparsity

- Regularization $\iff$ constraint:

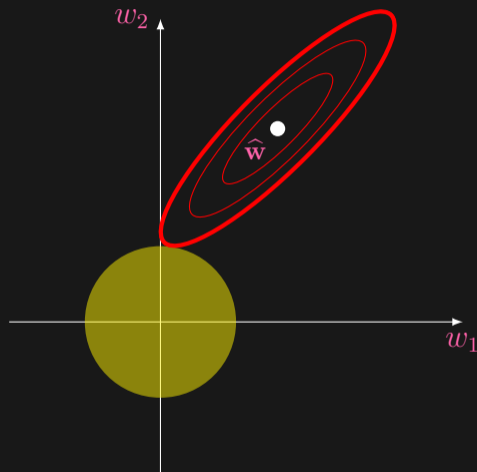$$\min_{\|\mathbf{W}\|_F \leq \gamma} \frac{1}{n} \|\mathbf{WX} - \mathbf{Y}\|_F^2$$

- Ridge regression $\to$ dense $\mathbf{W}$

- Lasso (Tibshirani, 1996):

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

- Regularization $\iff$ constraint:

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



R. Tibshirani. *"Regression Shrinkage and Selection via the Lasso"*. *Journal of the Royal Statistical Society: Series B*, vol. 58, no. 1 (1996), pp. 267–288.

# Sparsity

- Regularization $\Longleftrightarrow$ constraint:

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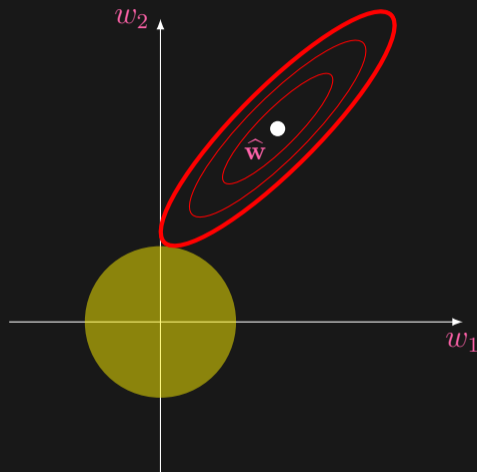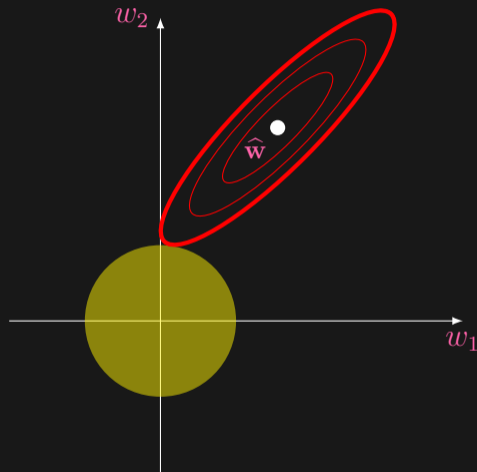
- Ridge regression $\to$ dense $\mathbf{W}$

- Lasso (Tibshirani, 1996):

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

- Regularization $\Longleftrightarrow$ constraint:

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



---

R. Tibshirani. "Regression Shrinkage and Selection via the Lasso". *Journal of the Royal Statistical Society: Series B*, vol. 58, no. 1 (1996), pp. 267–288.

# Sparsity

- Regularization $\Longleftrightarrow$ constraint:

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

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

- Lasso (Tibshirani, 1996):

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

- Regularization $\Longleftrightarrow$ constraint:

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



R. Tibshirani. "Regression Shrinkage and Selection via the Lasso". *Journal of the Royal Statistical Society: Series B*, vol. 58, no. 1 (1996), pp. 267–288.

# Sparsity

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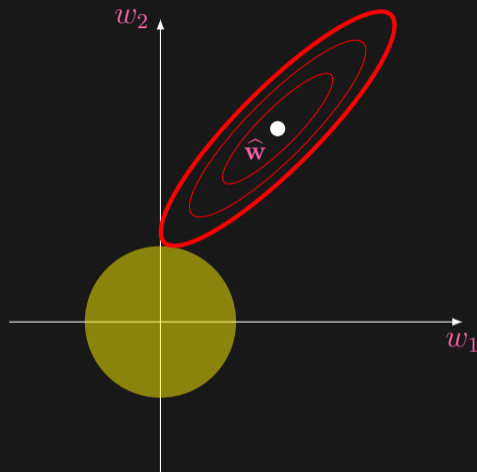
$$\min_{\|\mathbf{W}\|_F \leq \gamma} \frac{1}{n}\|\mathbf{WX} - \mathbf{Y}\|_F^2$$

- Ridge regression $\to$ 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:

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



R. Tibshirani. "Regression Shrinkage and Selection via the Lasso". *Journal of the Royal Statistical Society: Series B*, vol. 58, no. 1 (1996), pp. 267–288.
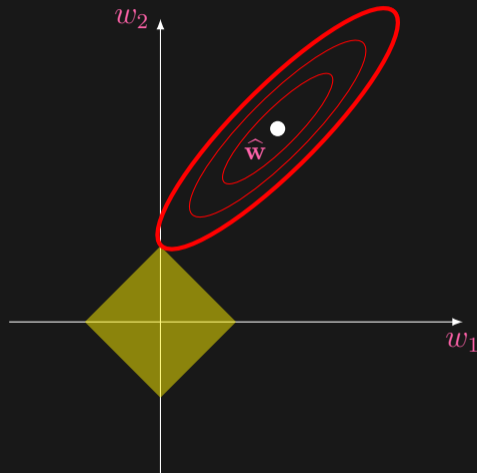
# Sparsity



- Regularization $\iff$ constraint:

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

- Ridge regression $\to$ 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:

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

R. Tibshirani. "Regression Shrinkage and Selection via the Lasso". *Journal of the Royal Statistical Society: Series B*, vol. 58, no. 1 (1996), pp. 267–288.
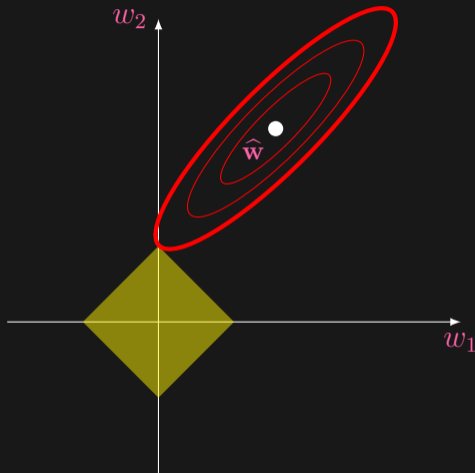
# Sparsity

- Regularization $\iff$ constraint:

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

- Ridge regression $\to$ 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:

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



R. Tibshirani. "Regression Shrinkage and Selection via the Lasso". *Journal of the Royal Statistical Society: Series B*, vol. 58, no. 1 (1996), pp. 267–288.

# Sparsity

- Regularization $\iff$ constraint:

$$\min_{\|\mathbf{W}\|_F \leq \gamma} \frac{1}{n}\|\mathbf{WX} - \mathbf{Y}\|_F^2$$
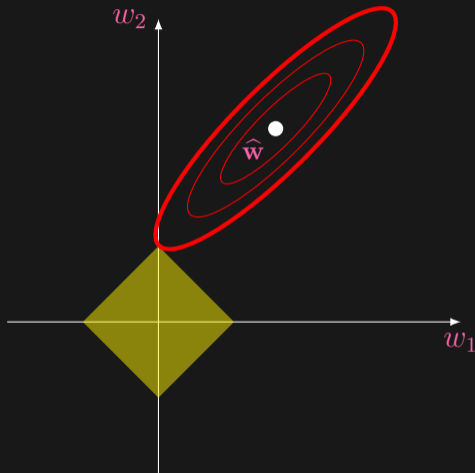
- Ridge regression $\rightarrow$ dense $\mathbf{W}$
  - more computation / communication
  - harder to interpret

- Lasso (Tibshirani, 1996):

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

- Regularization $\iff$ constraint:

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



R. Tibshirani. "Regression Shrinkage and Selection via the Lasso". *Journal of the Royal Statistical Society: Series B*, vol. 58, no. 1 (1996), pp. 267–288.

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

- 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 regularization:

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

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

R. Caruana. "Multitask Learning". *Machine Learning*, vol. 28 (1997), pp. 41–75, A. Argyriou et al. "Convex multi-task feature learning". *Machine Learning*, vol. 73 (2008), pp. 243–272.

# Task Regularization

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

- 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 regularization:

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

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

R. Caruana. "Multitask Learning". *Machine Learning*, vol. 28 (1997), pp. 41–75, A. Argyriou et al. "Convex multi-task feature learning". *Machine Learning*, vol. 73 (2008), pp. 243–272.
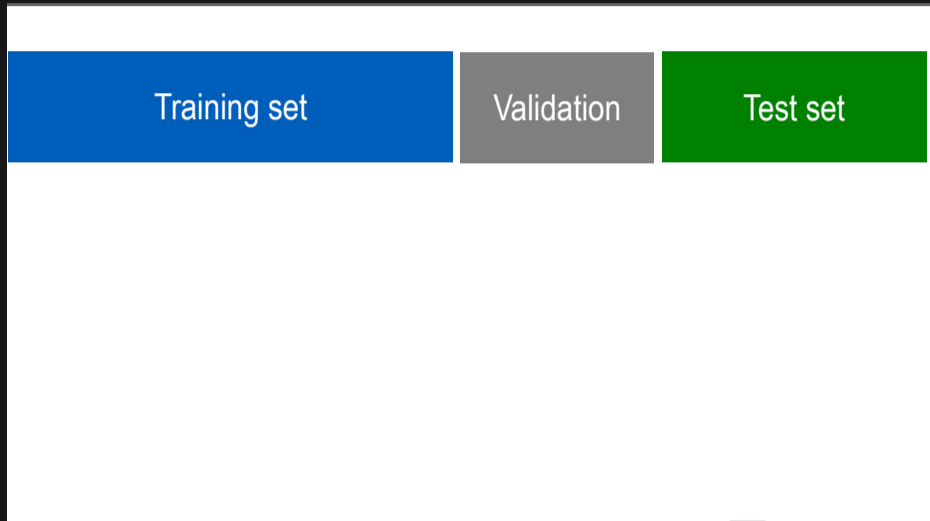
# Task Regularization

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

- In other words, the tasks are independent and can be solved separately

- Sometimes lumping tasks together (LHS) is computationally more efficient

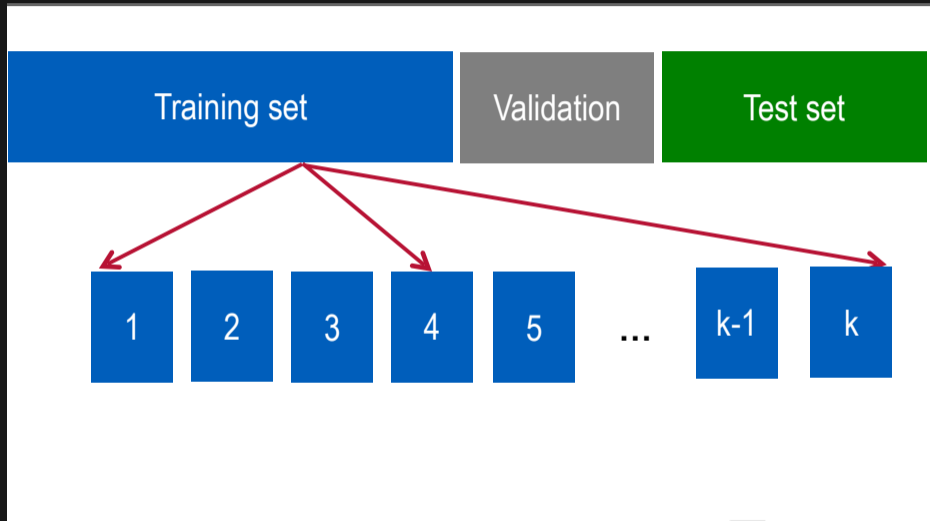- If tasks are related, can consider regularization:

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

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

R. Caruana. "Multitask Learning". *Machine Learning*, vol. 28 (1997), pp. 41–75, A. Argyriou et al. "Convex multi-task feature learning". *Machine Learning*, vol. 73 (2008), pp. 243–272.

# Task Regularization

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

- 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 regularization:

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

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

R. Caruana. "Multitask Learning". *Machine Learning*, vol. 28 (1997), pp. 41–75, A. Argyriou et al. "Convex multi-task feature learning". *Machine Learning*, vol. 73 (2008), pp. 243–272.
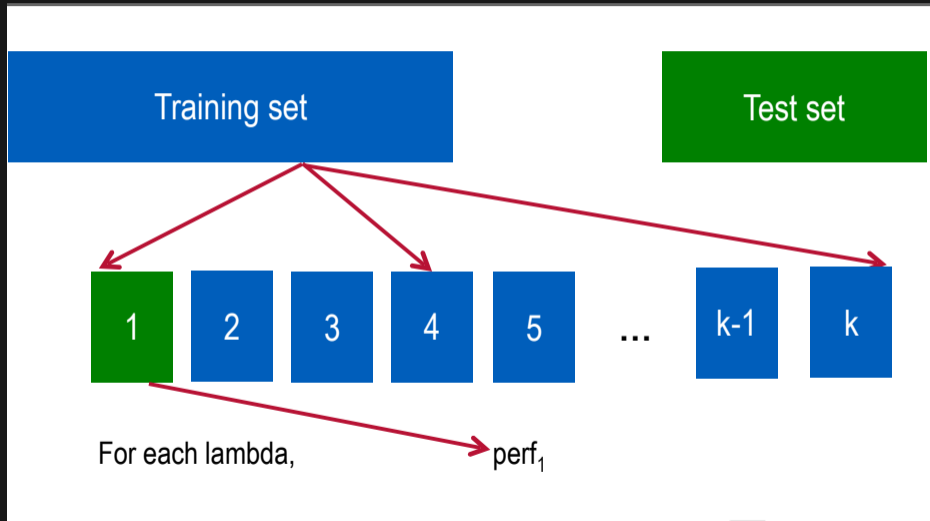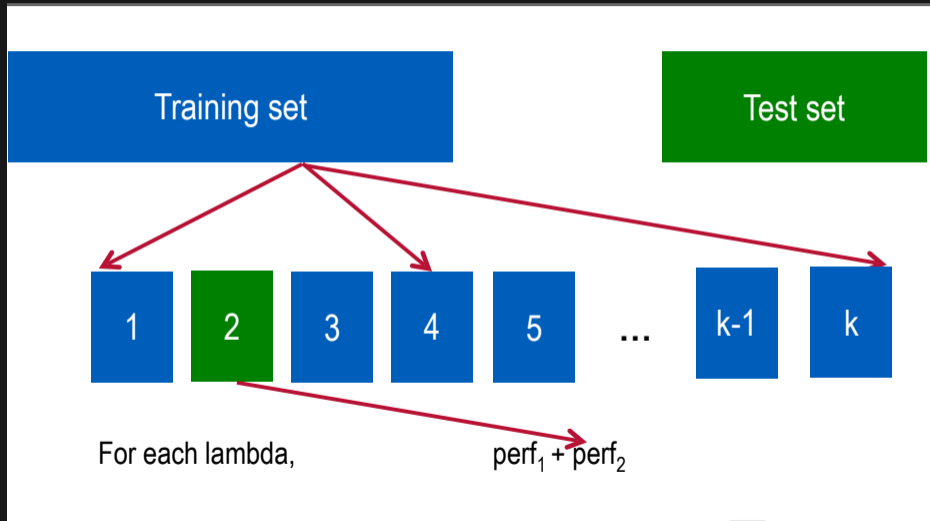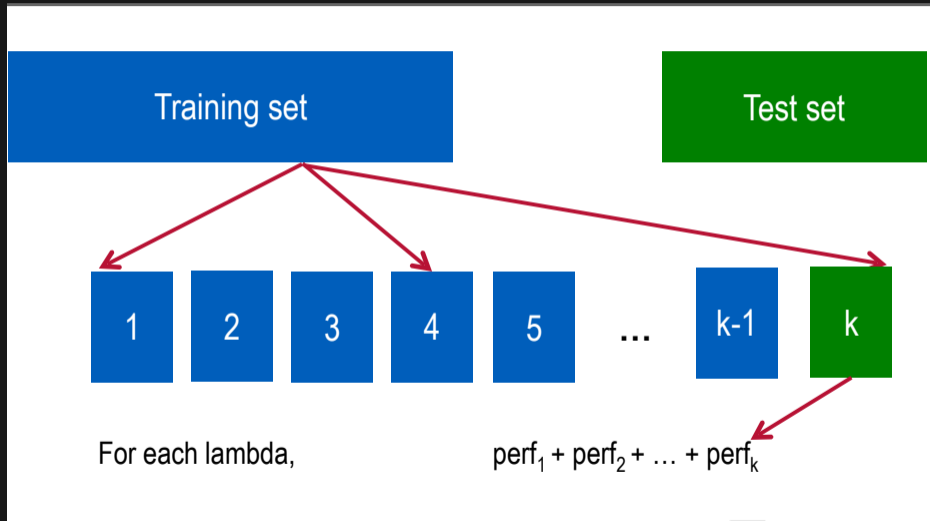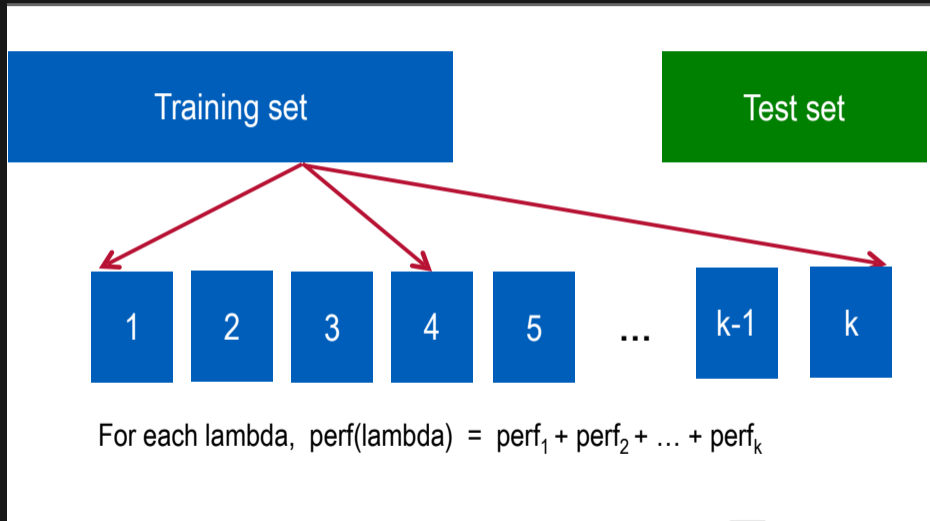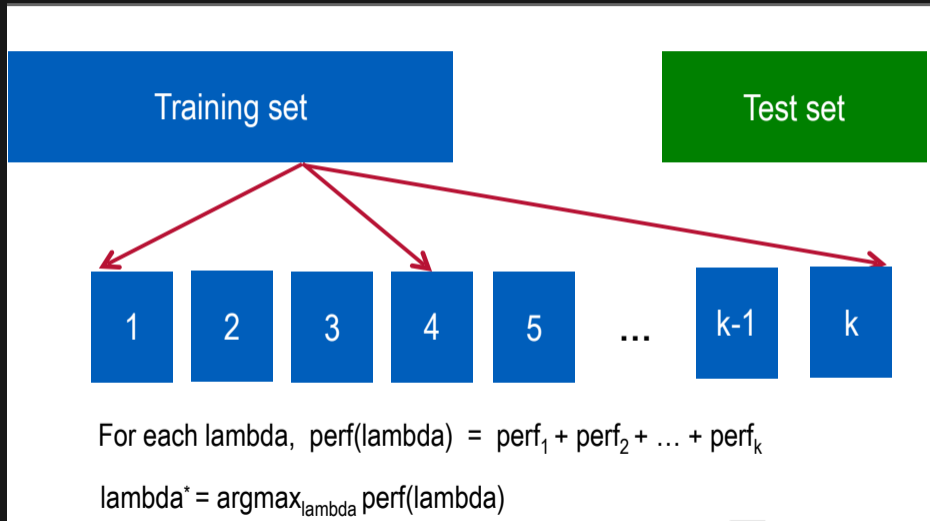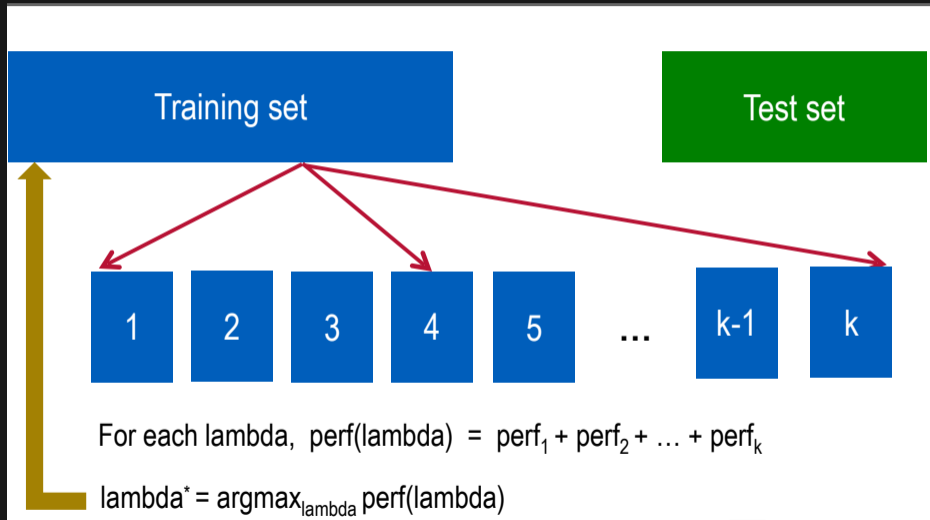
# Cross-validation

# Cross-validation

# Cross-validation

# Cross-validation

# Cross-validation



For each lambda, perf(lambda) = $perf_1 + perf_2 + \ldots + perf_k$

lambda* = $\text{argmax}_{lambda}$ perf(lambda)

For each lambda,  perf(lambda)  =  perf$_1$ + perf$_2$ + … + perf$_k$

lambda$^*$ = argmax$_{lambda}$ perf(lambda)

# Cross-validation



For each lambda, $\text{perf}(\text{lambda}) = \text{perf}_1 + \text{perf}_2 + \ldots + \text{perf}_k$

$\text{lambda}^* = \text{argmax}_{\text{lambda}} \, \text{perf}(\text{lambda})$