# Judging Test error

- Training error of a classifier f

$$\frac{1}{n} \sum_{i=1}^{n} 1_{f(X_i) \neq Y_i}$$

Training Data
$\{X_i, Y_i\}_{i=1}^{n}$

- What about test error?

    Can't compute it.

- How can we know classifier is not overfitting?

    Hold-out or Cross-validation

# Hold-out method

Can judge test error by using an independent sample of data.

## Hold – out procedure:

n data points available $\quad D \equiv \{X_i, Y_i\}_{i=1}^{n}$

1) Split into two sets (randomly and preserving label proportion):

Training dataset $\qquad$ Validation/Hold-out dataset

$$D_T = \{X_i, Y_i\}_{i=1}^{m} \qquad D_V = \{X_i, Y_i\}_{i=m+1}^{n}$$

often m = n/2

2) Train classifier on $D_T$. Report error on validation dataset $D_V$. Overfitting if validation error is much larger than training error

# Hold-out method

Drawbacks:

- May not have enough data to afford setting one subset aside for getting a sense of generalization abilities
- Validation error may be misleading (bad estimate of test error) if we get an "unfortunate" split

Limitations of hold-out can be overcome by a family of sub-sampling methods at the expense of more computation.
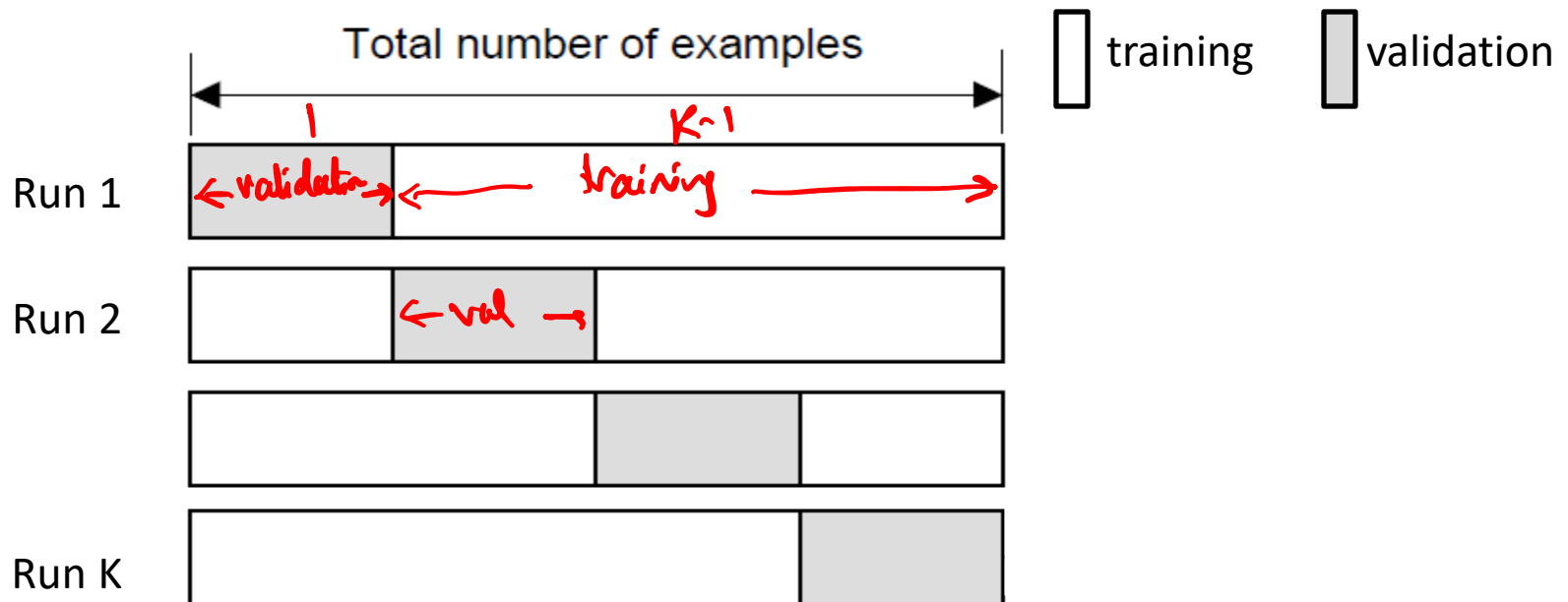
# Cross-validation

## K-fold cross-validation

Create K-fold partition of the dataset.
Do K runs: train using K-1 partitions and calculate validation error on remaining partition (rotating validation partition on each run).
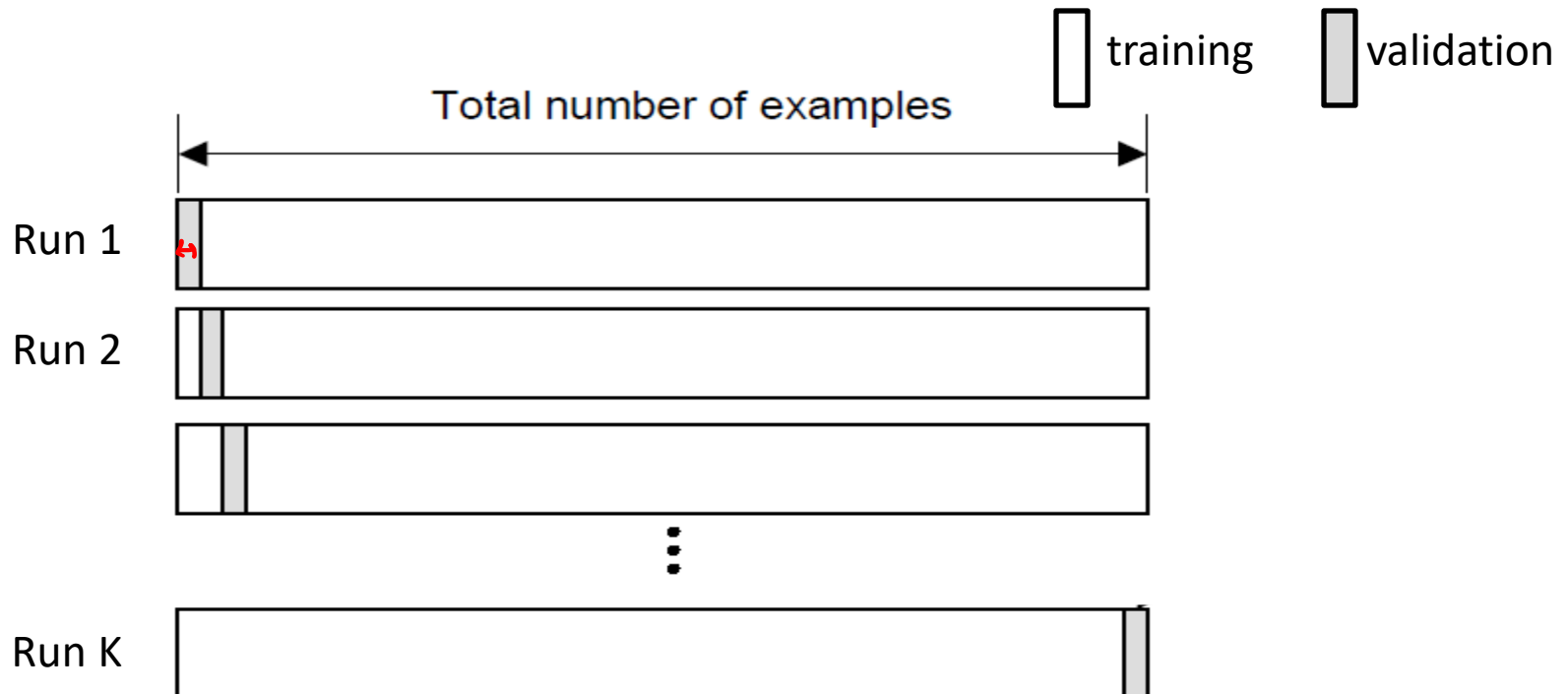Report average validation error

# Cross-validation

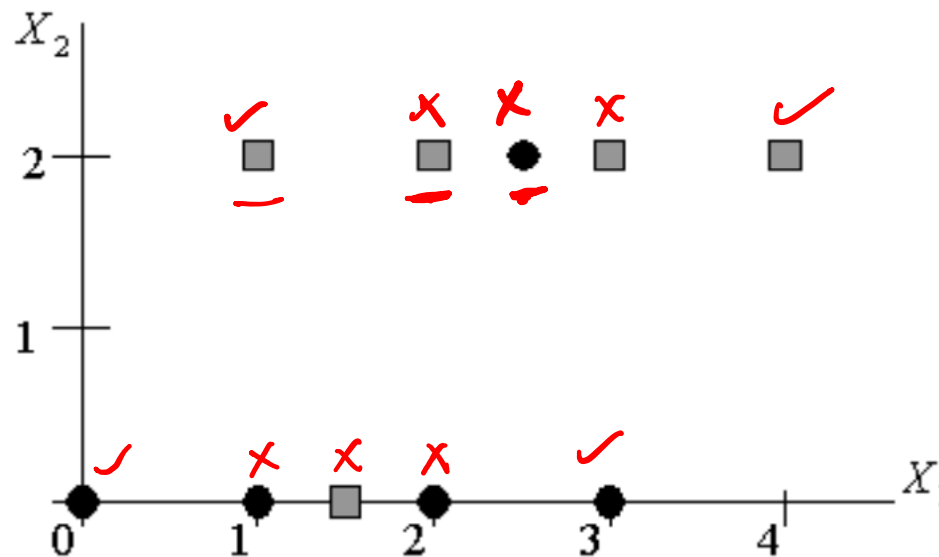## Leave-one-out (LOO) cross-validation

Special case of K-fold with K=n partitions
Equivalently, train on n-1 samples and validate on only one
sample per run for n runs

training   validation

Total number of examples

Run 1

Run 2

Run K

# Cross-validation

What is the leave-one-out cross-validation error of the given classifiers on the following dataset?



➢ Poll 1: Depth 1 Decision tree using best feature
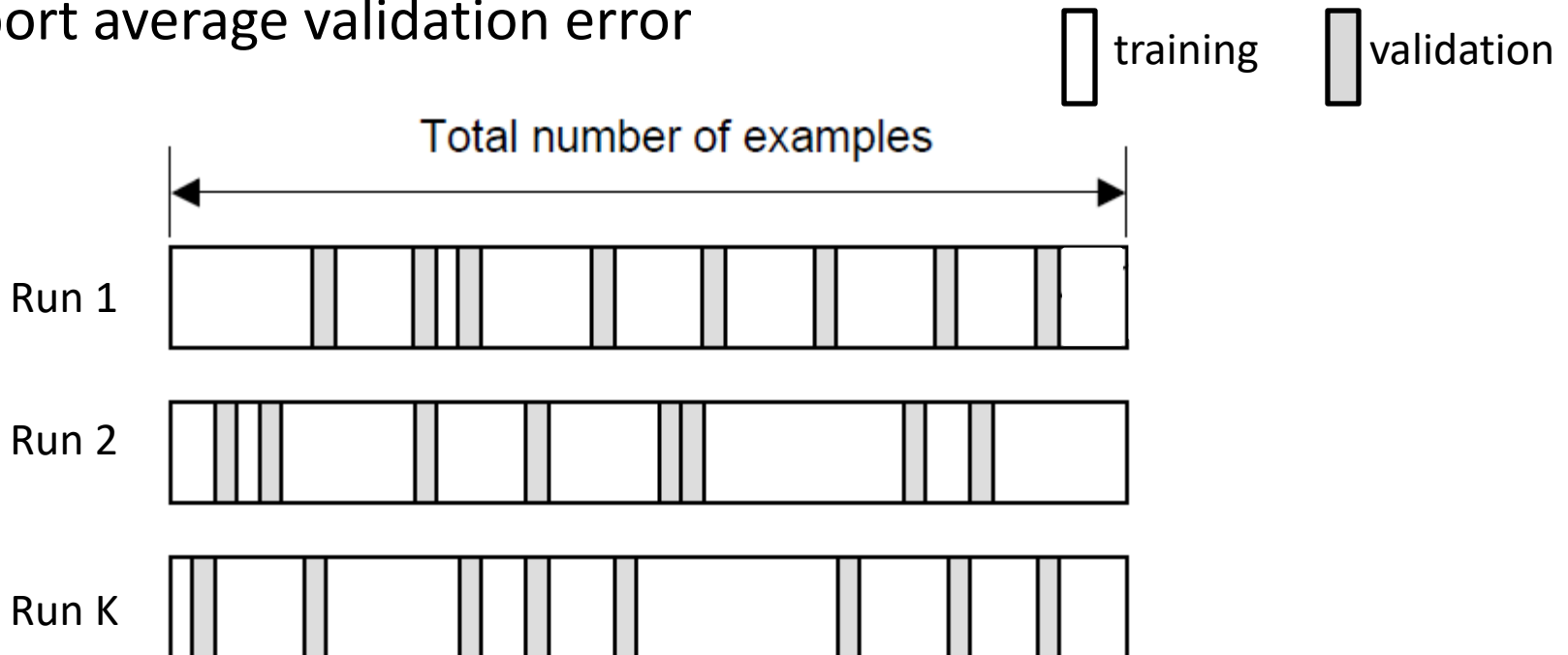➢ Poll 2: 1-NN classifier

# Cross-validation

## Random subsampling

Randomly subsample a fixed fraction $\alpha n$ (0< $\alpha$ <1) of the dataset for validation.

Compute validation error with remaining data as training data.

Repeat K times

Report average validation error

☐ training   ▨ validation



Total number of examples

Run 1

Run 2

Run K

# Practical Issues in Cross-validation

How to decide the values for *K* and *a*?

- Large K

    + Validation error can approximate test error well ✓

    - Observed validation error will be unstable (few validation pts)

    - The computational time will be very large as well (several runs)

- Small K

    + The #runs and, therefore, computation time are reduced

    + Observed validation error will be stable (many validation pts)

    - Validation error cannot approximate test error well ✓

Common choice: K = 10, $\alpha$ = 0.1 ☺

# Model selection using Hold-out/Cross-validation

- Train models of different complexities and evaluate their validation error using hold-out or cross-validation

- Pick model with smallest validation error (averaged over different runs for cross-validation)

➢ When using hold-out or cross-validation for model selection, test error should be reported using independent data

15

# What you should know

- Estimating test error using
  - hold-out
  - cross-validation
- Bias-variance tradeoff
- Model selection using
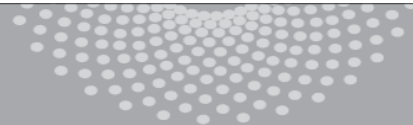  - hold-out
  - cross-validation

# Linear Regression

Aarti Singh

Machine Learning 10-701
Mar 20, 2023

# Supervised Learning Tasks

## Classification



X = Document      Y = Topic

Sports
Science
News

X = Cell Image      Y = Diagnosis

Anemic cell
Healthy cell

## Regression



Y = Age of a subject

X = Brain Scan

# Regression Tasks

Estimating Energy Usage

**X = building characteristics**
**Y = energy consumption**

energystar.gov

Estimating Contamination

**X = new location**
**Y = sensor reading**

# Mean Squared Error (MSE) Minimization

$|f(x) - Y|$

$= E[Y|x]$

Optimal predictor: $f^* = \arg\min_f \underset{XY}{\mathbb{E}}[\underbrace{(f(X) - Y)^2}]$

$$E[(f(x) - Y)^2] = E[(\underbrace{(f(x) - E[Y|x])}_{a} + \underbrace{(E[Y|x] - Y)}_{b})^2] \qquad (a+b)^2$$

$$= E[\underbrace{(f(x) - E[Y|x])^2}_{} + \underbrace{(E[Y|x] - Y)^2 + 2(f(x) - E[Y|x])}_{} $$
$$\underbrace{(E[Y|x] - Y)}_{c}]$$

$$c \quad E_{XY}[(f(x) - E[Y|x])(E[Y|x] - Y)]$$

$$= E_x[E_{Y|x}[(f(x) - E[Y|x])(E[Y|x] - Y)]] = 0$$

$$f^* = \arg\min_f E[(f(x) - E[Y|x])^2] = E[Y|x]$$

# Mean Squared Error (MSE) Minimization

Optimal predictor: $f^* = \arg\min_{f} \mathbb{E}[(f(X) - Y)^2] = \mathbb{E}[Y|X]$

Empirical Minimizer: $\widehat{f}_n = \arg\min_{f \in \mathcal{F}} \dfrac{1}{n} \sum_{i=1}^{n} (f(X_i) - Y_i)^2$

**Empirical mean**

Law of Large Numbers:

$$\dfrac{1}{n} \sum_{i=1}^{n} [\mathrm{loss}(Y_i, f(X_i))] \xrightarrow{\ n \longrightarrow \infty\ } \mathbb{E}_{XY}[\mathrm{loss}(Y, f(X))]$$

# Restrict class of predictors

Optimal predictor:

$$f^* = \arg\min_f \mathbb{E}[(f(X) - Y)^2]$$

Empirical Minimizer:

$$\widehat{f}_n = \arg\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} (f(X_i) - Y_i)^2$$

**Class of predictors**

➢ Why?

$\mathcal{F}$ - Class of Linear functions
  - Class of Polynomial functions
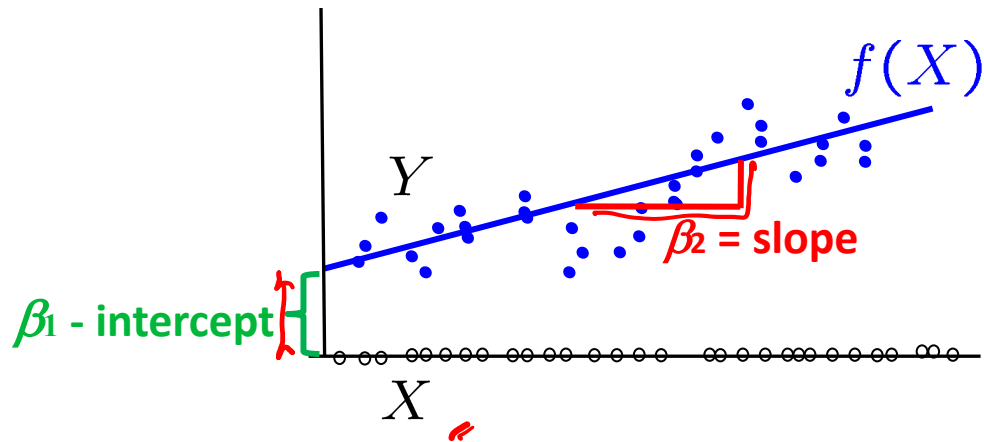  - Class of nonlinear functions

# Linear Regression

$$\widehat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^{n} (f(X_i) - Y_i)^2$$

Least Squares Estimator

$\mathcal{F}_L$ - Class of Linear functions

Uni-variate case:

$$f(X) = \beta_1 + \beta_2 X$$

$f(X)$

$Y$

$\beta_2$ = slope

$\beta_1$ - intercept

$X$

Multi-variate case:

p features

$$f(X) = f(X^{(1)}, \dots, X^{(p)}) = \beta_1 X^{(1)} + \beta_2 X^{(2)} + \dots + \beta_p X^{(p)}$$

1

$$= X\beta \qquad \text{where} \qquad X = [X^{(1)} \dots X^{(p)}], \quad \beta = [\beta_1 \dots \beta_p]^T$$

# Linear Regression (Matrix-vector form)

$$\widehat{f}_n^L = \arg\min_{f \in \mathcal{F}_L} \frac{1}{n}\sum_{i=1}^n (f(X_i) - Y_i)^2 \qquad\qquad f(X_i) = X_i\beta$$

$$\widehat{\beta} = \arg\min_\beta \frac{1}{n}\sum_{i=1}^n (X_i\beta - Y_i)^2 \qquad\qquad \widehat{f}_n^L(X) = X\widehat{\beta}$$

$$= \arg\min_\beta \frac{1}{n}(\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y})$$

$$A\beta = \begin{bmatrix} X_1\beta \\ X_2\beta \\ \dot{x}_n\beta \end{bmatrix}_{n\times 1}$$

$$\text{training data matrix} \quad \mathbf{A} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \ldots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \ldots & X_n^{(p)} \end{bmatrix} \quad \beta \quad \mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_n \end{bmatrix}$$

$n \times p$

$n \times 1$

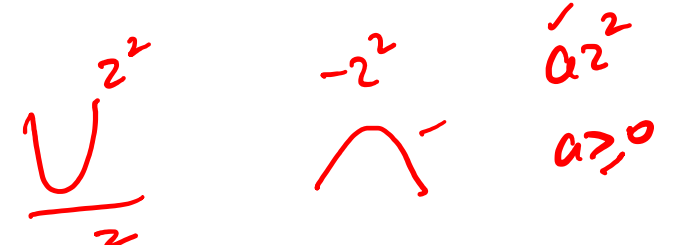$\underset{n\times p}{A}\ \underset{p\times 1}{\beta}$

# Linear Regression

$$\widehat{\beta} = \arg\min_{\beta} \frac{1}{n}(\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) = \arg\min_{\beta} J(\beta)$$

$$J(\beta) = (\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y})$$

$$= (\beta^T A^T - Y^T)(A\beta - Y)$$

$$= \beta^T A^T A \beta \cdots$$

$A =$ training data matrix

➤ Poll

Is the objective convex in β?

A) Convex, quadratic in β
B) Non-convex, A may not be positive semi definite
C) Depends on conditioning (ratio of max:min eigenvalues) of $A^TA$
D) Convex, $A^TA$ is positive semi definite

$$v^T A^T A v = \|Av\|^2 \qquad \frac{v^T M v}{v^T v}$$

# Linear Regression

$$① \quad \frac{\partial z^T \beta}{\partial \beta} = z \qquad ② \quad \frac{\partial \beta^T M \beta}{\partial \beta} = \underbrace{(M + M^T)}_{=2M \text{ if } M = M^T} \beta$$

$$\widehat{\beta} = \arg\min_\beta \frac{1}{n}(\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) = \arg\min_\beta J(\beta)$$

$$J(\beta) = (\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) = \beta^T \underline{A^T A} \beta - \underbrace{Y^T A}_{z^T} \beta - \underbrace{\beta^T A^T Y}_{= Y^T A \beta} + Y^T Y$$

$$\frac{\partial J(\beta)}{\partial \beta} = 2(A^T A)\beta - 2A^T Y \Big|_{\widehat{\beta}} = 0$$

$$\frac{\partial J(\beta)}{\partial \beta}\Big|_{\widehat{\beta}} = 0 \qquad (A^T A)\widehat{\beta} = A^T Y$$

# Linear regression solution satisfies Normal Equations

$$(\mathbf{A}^T\mathbf{A})\widehat{\beta} = \mathbf{A}^T\mathbf{Y}$$

$Y = \begin{bmatrix} Y_1 \\ \vdots \\ y_n \end{bmatrix}$  $A = \begin{bmatrix} x_1^{(1)} \cdots x_1^{(p)} \\ \\ x_n^{(1)} \cdots x_n^{(p)} \end{bmatrix}$
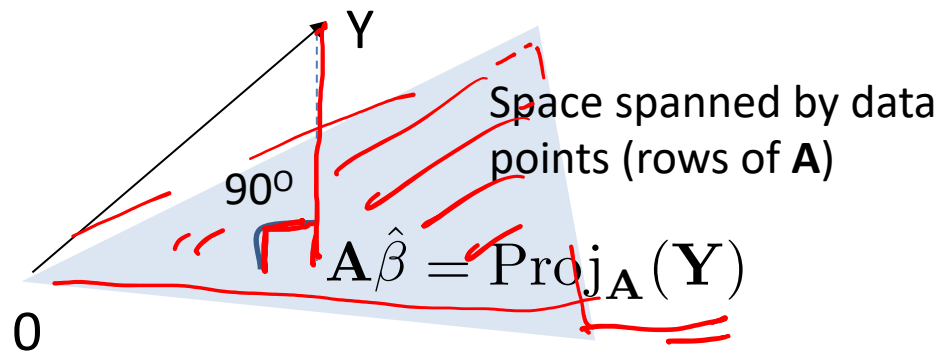
p x p  p x1        p x1

If $(\mathbf{A}^T\mathbf{A})$ is invertible,

$$\widehat{\beta} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{Y} \qquad \widehat{f}_n^L(X) = X\widehat{\beta}$$

$\widehat{f}_n^L(A) = A\widehat{\beta} = A(A^TA)A^TY$
$= \text{Proj}_A(Y)$

Predicted labels for training points $\mathbf{A}\widehat{\beta} = \text{Proj}_{\mathbf{A}}(\mathbf{Y})$



Y

Space spanned by data points (rows of **A**)

90º

$\mathbf{A}\widehat{\beta} = \text{Proj}_{\mathbf{A}}(\mathbf{Y})$

0

# Linear regression solution satisfies Normal Equations

$$(\mathbf{A}^T\mathbf{A})\widehat{\beta} = \mathbf{A}^T\mathbf{Y}$$

p x p   p x1       p x1

If $(\mathbf{A}^T\mathbf{A})$ is invertible,

$$\widehat{\beta} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{Y} \qquad \widehat{f}_n^L(X) = X\widehat{\beta}$$

Later: When is $(\mathbf{A}^T\mathbf{A})$ invertible ?

Now: What if $(\mathbf{A}^T\mathbf{A})$ is invertible but expensive (p very large)?

# Gradient Descent

Even when $(\mathbf{A}^T \mathbf{A})$ is invertible, might be computationally expensive if **A** is huge.

$$\widehat{\beta} = \arg\min_{\beta} \frac{1}{n}(\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) = \arg\min_{\beta} J(\beta)$$
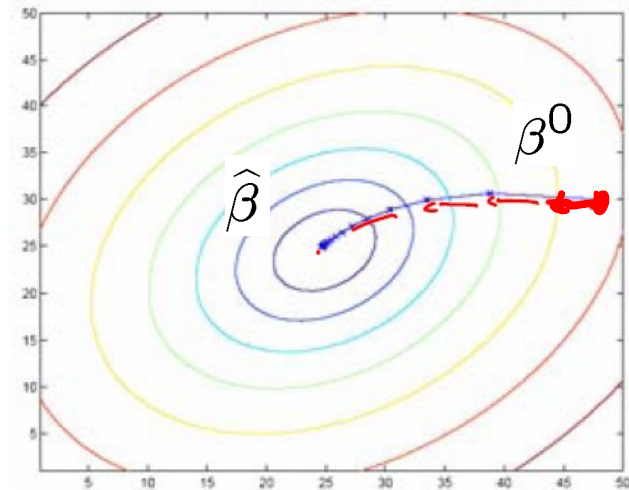
**Since J(β) is convex, move along negative of gradient**

Initialize: $\beta^0$

step size

Update: $\beta^{t+1} = \beta^t - \dfrac{\alpha}{2}\dfrac{\partial J(\beta)}{\partial \beta}\Big|_t$

predicted labels

$= \beta^t - \alpha\, \mathbf{A}^T(\mathbf{A}\beta^t - Y)$

0 if $\widehat{\beta} = \beta^t$

$\beta^0$

$\widehat{\beta}$

Stop: when some criterion met e.g. fixed # iterations, or $\dfrac{\partial J(\beta)}{\partial \beta}\Big|_{\beta^t} < \varepsilon$.

13

# Least Square solution satisfies Normal Equations

$$\left.\frac{\partial J(\beta)}{\partial \beta}\right|_{\widehat{\beta}} = 0 \qquad \text{gives} \qquad (\mathbf{A}^T\mathbf{A})\widehat{\beta} = \mathbf{A}^T\mathbf{Y}$$

<span style="color:blue">p x p   p x1     p x1</span>

If $(\mathbf{A}^T\mathbf{A})$ is invertible,

1) If dimension p not too large, analytical solution:

$$\widehat{\beta} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{Y} \qquad \widehat{f}_n^L(X) = X\widehat{\beta}$$

2)  If dimension p is large, computing inverse is expensive O(p³)
    Gradient descent since objective is convex ($\mathbf{A}^\mathsf{T}\mathbf{A} \succeq 0$)

$$\begin{aligned} \beta^{t+1} &= \beta^t - \frac{\alpha}{2}\left.\frac{\partial J(\beta)}{\partial \beta}\right|_t \\ &= \beta^t - \alpha\, \mathbf{A}^T(\mathbf{A}\beta^t - Y) \end{aligned}$$