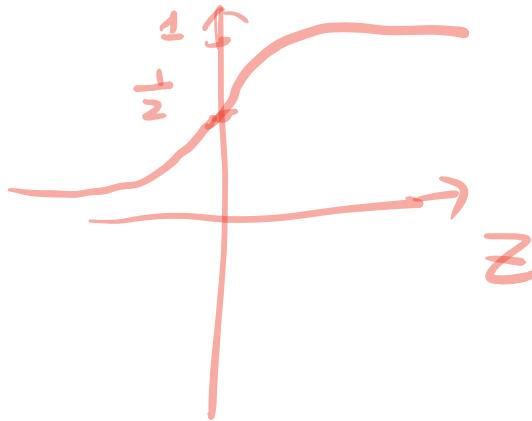# Linear Regression

# Announcements

# Recap on Logistic Regression / Optimization

Binary classification with $\mathcal{D} = \{x_i, y_i\}_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}$

Logistic Regression asumes $P(y \mid x; w) = \dfrac{1}{1 + \exp(-y(w^\top x))}$ ✓

$z := y(w^\top x)$

# Recap on Logistic Regression / Optimization

Binary classification with $\mathcal{D} = \{x_i, y_i\}_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in \{-1,1\}$

Logistic Regression asumes $P(y \mid x; w) = \dfrac{1}{1 + \exp(-y(w^\top x))}$

<span style="color:red">Using MLE, we get our optimization objective:</span>

$$\hat{w} := \arg\min_w \sum_{i=1}^n \ln(1 + \exp(-y_i(w^\top x_i)))$$

# Recap on Logistic Regression / Optimization

Binary classification with $\mathcal{D} = \{x_i, y_i\}_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in \{-1,1\}$

Logistic Regression asumes $P(y \,|\, x; w) = \dfrac{1}{1 + \exp(-y(w^\top x))}$

Using MLE, we get our optimization objective:

$$\hat{w} := \arg\min_{w} \sum_{i=1}^{n} \ln(1 + \exp(-y_i(w^\top x_i)))$$

Given a test example $x_{test}$, we can make prediction:

$$\hat{y} = \begin{cases} +1 & P(+1 \,|\, x_{test}; \hat{w}) > 0.5 \\ -1 & \text{else} \end{cases}$$

# Recap on Logistic Regression / Optimization

$$\arg\min_{w} \sum_{i=1}^{n} \ln(1 + \exp(-y_i(w^\top x_i)))$$

**Logistic Regression with SGD as the optimizer:**

Initialize $w^0 = \mathbf{0}$

While not converged:

# Recap on Logistic Regression / Optimization

$$\arg\min_{w} \sum_{i=1}^{n} \ln(1 + \exp(-y_i(w^\top x_i)))$$

**Logistic Regression with SGD as the optimizer:**

Initialize $w^0 = \mathbf{0}$

While not converged:

  Randomly sample one training pair $(x, y) \sim \mathcal{D}$

$$\nabla_w \ln\left(1 + \exp(-y(w^\top x))\right)$$

# Recap on Logistic Regression / Optimization

$$\arg \min_{w} \sum_{i=1}^{n} \ln(1 + \exp(-y_i(w^\top x_i)))$$

**Logistic Regression with SGD as the optimizer:**

Initialize $w^0 = \mathbf{0}$

While not converged:

Randomly sample one training pair $(x, y) \sim \mathcal{D}$

Compute gradient: $g = (1 - P(y \mid x; w^t)) \cdot (-yx)$

# Recap on Logistic Regression / Optimization

$$\arg \min_{w} \sum_{i=1}^{n} \ln(1 + \exp(-y_i(w^\top x_i)))$$

**Logistic Regression with SGD as the optimizer:**

Initialize $w^0 = \mathbf{0}$

While not converged:

> Randomly sample one training pair $(x, y) \sim \mathcal{D}$
>
> Compute gradient: $g = (1 - P(y \mid x; w^t)) \cdot (-yx)$
>
> SGD update: $w^{t+1} = w^t + \alpha(1 - P(y \mid x; w^t)) \cdot (yx)$

# Recap on Logistic Regression / Optimization

$$\arg\min_{w} \sum_{i=1}^{n} \ln(1 + \exp(-y_i(w^\top x_i)))$$

**Logistic Regression with SGD as the optimizer:**

Initialize $w^0 = \mathbf{0}$

While not converged:

    Randomly sample one training pair $(x, y) \sim \mathcal{D}$

    Compute gradient: $g = (1 - P(y \mid x; w^t)) \cdot (-yx)$

    SGD update: $w^{t+1} = w^t + \alpha(1 - P(y \mid x; w^t)) \cdot (yx)$

Compare this to Perceptron!

# Outline for Today

1. Intro on Linear Regression

2. Normal equation for linear Regression

3. Interpretation of Linear Regression using MLE / MAP
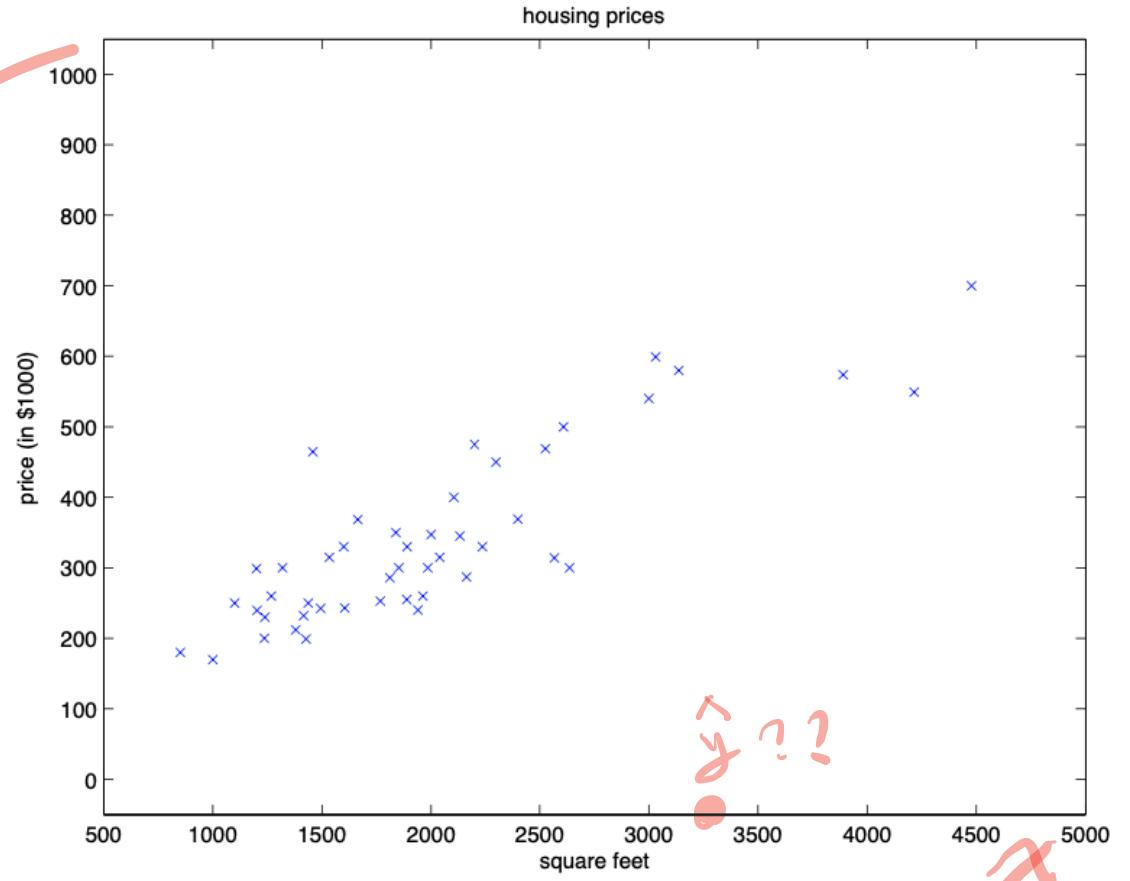
# Ex: Predicting the house price

Dataset:

| Living area $(\text{feet}^2)$ | Price $(1000\$\text{s})$ |
|:---:|:---:|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| $\vdots$ | $\vdots$ |
| $x$ | $y$ |

(Example from Stanford CS229)

# Ex: Predicting the house price

### Dataset:

| Living area (feet$^2$) | Price (1000$s) |
|:---:|:---:|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| ⋮ | ⋮ |
| $x$ | $y$ |

(Example from Stanford CS229)

### Plot:



housing prices

# Ex: Predicting the house price

### Dataset:

| Living area ($\text{feet}^2$) | Price (1000\$s) |
|:---:|:---:|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| ⋮ | ⋮ |
| $x$ | $y$ |

(Example from Stanford CS229)

### Plot:



housing prices

$$\ell = \left(\vec{w}^\top x - y\right)^2$$

$$h(x) = w_1 x + w_0$$

# Ex: Predicting the house price (2d case)

Dataset:

| Living area (feet$^2$) | #bedrooms | Price (1000\$s) |
|:---:|:---:|:---:|
| 2104 | 3 | 400 |
| 1600 | 3 | 330 |
| 2400 | 3 | 369 |
| 1416 | 2 | 232 |
| 3000 | 4 | 540 |
| $\vdots$ | $\vdots$ | $\vdots$ |
| $x[1]$ | $x[2]$ | $y$ |

(Example from Stanford CS229)

# Ex: Predicting the house price (2d case)

Dataset:

| Living area (feet$^2$) | #bedrooms | Price (1000$s) |
|:---:|:---:|:---:|
| 2104 | 3 | 400 |
| 1600 | 3 | 330 |
| 2400 | 3 | 369 |
| 1416 | 2 | 232 |
| 3000 | 4 | 540 |
| $\vdots$ | $\vdots$ | $\vdots$ |
| $x[1]$ | $x[2]$ | $y$ |

Goal: finding the linear function

$$h(x) = w_1 x[1] + w_2 x[2] + w_0$$

that fits the data well

(Example from Stanford CS229)

# Ex: Predicting the house price (2d case)

Dataset:

| Living area (feet$^2$) | #bedrooms | Price (1000$s) |
|:---:|:---:|:---:|
| 2104 | 3 | 400 |
| 1600 | 3 | 330 |
| 2400 | 3 | 369 |
| 1416 | 2 | 232 |
| 3000 | 4 | 540 |
| $\vdots$ | $\vdots$ | $\vdots$ |
| $x[1]$ | $x[2]$ | $y$ |

As usual, we append 1 to the feature, i.e.,

$$x = \begin{bmatrix} x[1] \\ x[2] \\ 1 \end{bmatrix}$$

So the linear function can be written as:

$$h(x) = w^\top x$$

$$= \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}^\top \begin{bmatrix} x[1] \\ x[2] \\ 1 \end{bmatrix}$$

# Outline for Today

1. Intro on Linear Regression

2. Normal equation for linear Regression

3. Interpretation of Linear Regression using MLE / MAP

# Mathematical formulation of linear regression

**Input**: dataset $\mathscr{D} = \{x_i, y_i\}_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

# Mathematical formulation of linear regression

**Input**: dataset $\mathscr{D} = \{x_i, y_i\}_{i=1}^{n}, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

**Hypothesis**: linear function $h(x) = w^\top x$

$w \in \mathbb{R}^d$

# Mathematical formulation of linear regression

**Input**: dataset $\mathscr{D} = \{x_i, y_i\}_{i=1}^{n}, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

**Hypothesis**: linear function $h(x) = w^\top x$

**Hypothesis class**: all possible linear functions $\{w^\top x, \forall w \in \mathbb{R}^d\}$

# Mathematical formulation of linear regression

**Input**: dataset $\mathcal{D} = \{x_i, y_i\}_{i=1}^{n}, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

**Hypothesis**: linear function $h(x) = w^\top x$

**Hypothesis class**: all possible linear functions $\{w^\top x, \forall w \in \mathbb{R}^d\}$

**Loss function**: squared loss $\ell(w^\top x, y) = (w^\top x - y)^2$

Diff between $w^\top x$ & GT $y$

# Mathematical formulation of linear regression

**Input**: dataset $\mathscr{D} = \{x_i, y_i\}_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

**Hypothesis**: linear function $h(x) = w^\top x$

**Hypothesis class**: all possible linear functions $\{w^\top x, \forall w \in \mathbb{R}^d\}$

**Loss function**: squared loss $\ell(w^\top x, y) = (w^\top x - y)^2$

Q: can we use absolute loss, i.e., $|w^\top x - y|$ ?

# Mathematical formulation of linear regression

**Input**: dataset $\mathcal{D} = \{x_i, y_i\}_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

**Hypothesis**: linear function $h(x) = w^\top x$

**Hypothesis class**: all possible linear functions $\{w^\top x, \forall w \in \mathbb{R}^d\}$

**Loss function**: squared loss $\ell(w^\top x, y) = (w^\top x - y)^2$

Q: can we use absolute loss, i.e., $|w^\top x - y|$ ?

Q: can we use $(w^\top x - y)^3$ ?

$\left(w^\top x - y\right)^{10}$

# Mathematical formulation of linear regression

Formulating the optimization problem:

$$\arg\min_{w} \sum_{i=1}^{n} (w^{\top} x_i - y_i)^2$$

# Linear regression solution

$$\arg\min_w \sum_{i=1}^{n} (w^\top x_i - y_i)^2$$

Let's compute the closed-form solution:

# Linear regression solution

$$\arg\min_{w} \sum_{i=1}^{n} (w^\top x_i - y_i)^2$$

Let's compute the closed-form solution:

Define $X = [x_1, x_2, \ldots, x_n] \in \mathbb{R}^{d \times n}$, $Y = [y_1, \ldots, y_n]^\top \in \mathbb{R}^n$

$$X = \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \cdots & x_n \\ | & | & & | \end{bmatrix} \qquad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

# Linear regression solution

$$\arg\min_w \sum_{i=1}^{n} (w^\top x_i - y_i)^2$$

$$X = \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \cdots & x_n \\ | & | & & | \end{bmatrix}$$

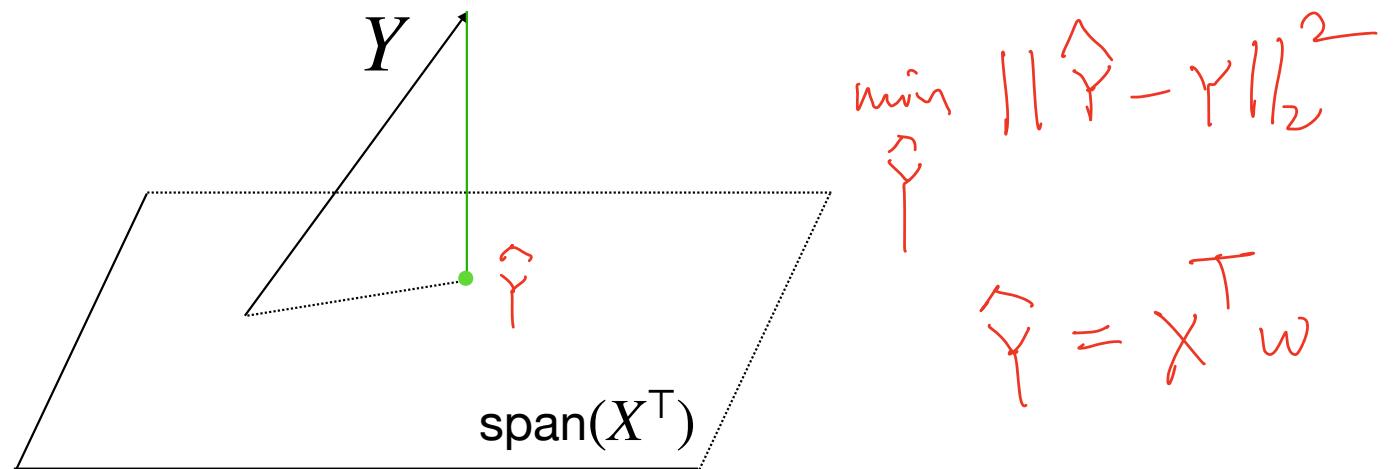$$\|w\|_2^2 = \sum_{i=1}^{d} w_i^2$$

Let's compute the closed-form solution:

Define $X = [x_1, x_2, \ldots, x_n] \in \mathbb{R}^{d \times n}$, $Y = [y_1, \ldots, y_n]^\top \in \mathbb{R}^n$

$$\sum_{i=1}^{n} (w^\top x_i - y_i)^2 = \|X^\top w - Y\|_2^2$$

$$\left\| \begin{bmatrix} - & x_1^\top & - \\ - & x_n^\top & - \\ & \vdots & \\ - & x_n^\top & - \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \right\|_2^2$$

# Linear regression solution

$$\arg\min_{w} \sum_{i=1}^{n} (w^\top x_i - y_i)^2$$

Let's compute the closed-form solution:

Define $X = [x_1, x_2, \ldots, x_n] \in \mathbb{R}^{d \times n}$, $Y = [y_1, \ldots, y_n]^\top \in \mathbb{R}^n$

$$\sum_{i=1}^{n} (w^\top x_i - y_i)^2 = \|X^\top w - Y\|_2^2$$

$$\Rightarrow \arg\min_{w} \|X^\top w - Y\|_2^2$$

# Linear regression solution

$$\arg \min_{w} \|X^{\top}w - Y\|_2^2 \qquad X = [x_1, x_2, \ldots, x_n] \in \mathbb{R}^{d \times n}$$

$Y$

$\hat{Y}$

span$(X^{\top})$

$$\min_{\hat{Y}} \|\hat{Y} - Y\|_2^2$$

$$\hat{Y} = X^{\top}w$$

# Linear regression solution

$$\arg \min_{w} \|X^\top w - Y\|_2^2$$

# Linear regression solution

$$\left(X^\top w - Y\right)^\top \left(X^\top w - Y\right)$$

$$\arg \min_{w} \|X^\top w - Y\|_2^2$$

$$\nabla_w \|X^\top w - Y\|_2^2 = XX^\top w - XY$$

$$2 \qquad 2 = 0$$

$$\| X^\top w - Y \|_2^2 = w^\top X X^\top w - 2 w^\top X Y + Y^\top Y$$

$$\nabla_w \left( w^\top X X^\top w \right) = X X^\top w + \left( X X^\top \right)^\top w$$

$$= 2 \left( X X^\top \right) w$$

$$\nabla_w \left( 2 w^\top X Y \right) = 2 X Y$$

# Linear regression solution

$$\arg \min_{w} \|X^\top w - Y\|_2^2$$

$$\nabla_w \|X^\top w - Y\|_2^2 = XX^\top w - XY = 0 \qquad XX^\top w = XY$$

$$\underline{\qquad\qquad\qquad\qquad}$$

Normal Eqn

if $XX^\top$ is full rank, then $\hat{w} = (XX^\top)^{-1}XY$

$$\underbrace{(XX^\top)^{-1}XX^\top}_{= I} w = (XX^\top)^{-1}XY$$

$$\Rightarrow w = (XX^\top)^{-1}XY$$

# Linear regression solution

$$\arg\min_{w} \|X^\top w - Y\|_2^2$$

$$\nabla_w \|X^\top w - Y\|_2^2 = XX^\top w - XY$$

if $XX^\top$ is full rank, then $\hat{w} = (XX^\top)^{-1}XY$

$Y$

$\|\hat{Y} - Y\|_2^2$

$\hat{Y} = X^\top \hat{w} = X^\top (XX^\top)^{-1} XY$

$\hat{Y} = X^\top \hat{w}$

$\text{span}(X^\top)$

# Linear regression solution

$$\arg \min_{w} \|X^\top w - Y\|_2^2$$

$$\nabla_w \|X^\top w - Y\|_2^2 = XX^\top w - XY$$

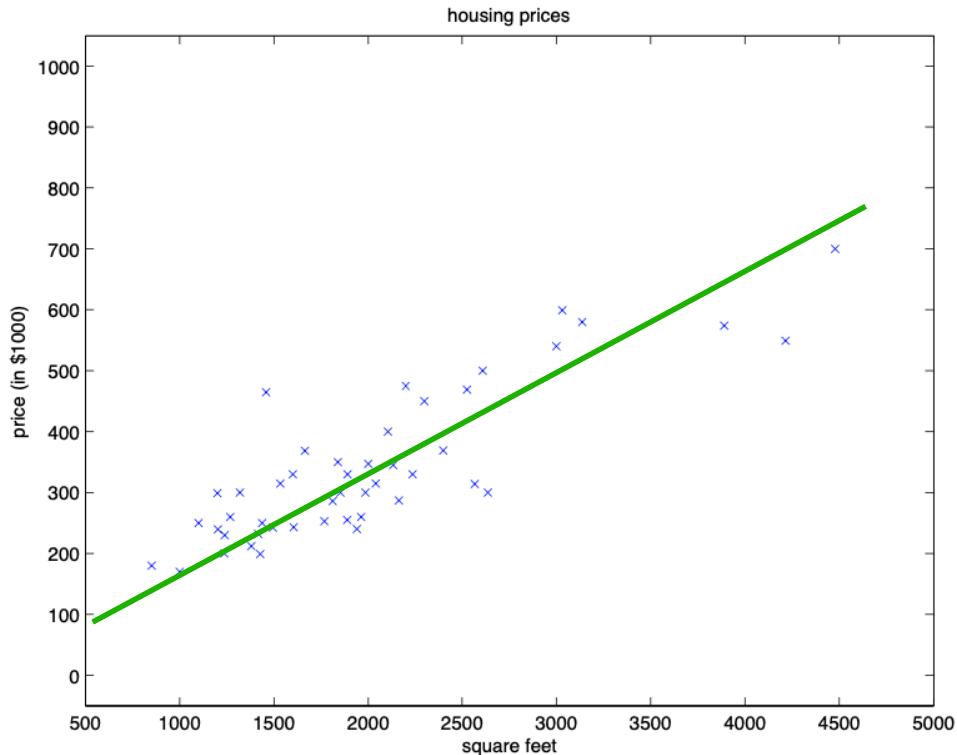if $XX^\top$ is full rank, then $\hat{w} = (XX^\top)^{-1} XY$



$Y$

$\hat{Y} = X^\top \hat{w}$

$\text{span}(X^\top)$

$X = \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \cdots & x_u \\ | & | & & | \end{bmatrix}$

$\in \mathbb{R}^{d \times n}$

What if $XX^\top$ is not full rank?

$XX^\top \in \mathbb{R}^{d \times d}$

# Linear regression solution

$$\arg \min_w \|X^\top w - Y\|_2^2$$

$$\nabla_w \|X^\top w - Y\|_2^2 = XX^\top w - XY$$

if $XX^\top$ is full rank, then $\hat{w} = (XX^\top)^{-1}XY$

$Y$

$\hat{Y} = X^\top \hat{w}$

$\text{span}(X^\top)$

What if $XX^\top$ is not full rank?

(We will talk about regularization soon)

# Prediction using linear regression

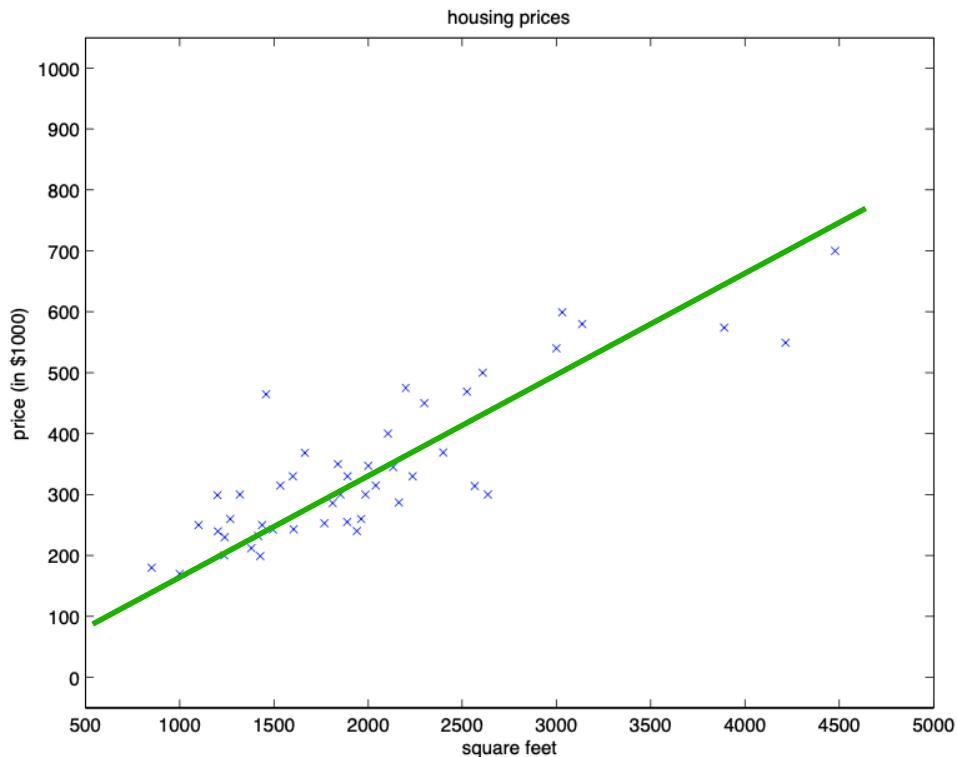Once we learned $\hat{w}$, we can use it to make prediction on any new feature $x$



Given $x_{test}$, our prediction is:

$$\hat{y} = x_{test}^{\top} \hat{w}$$

# Prediction using linear regression

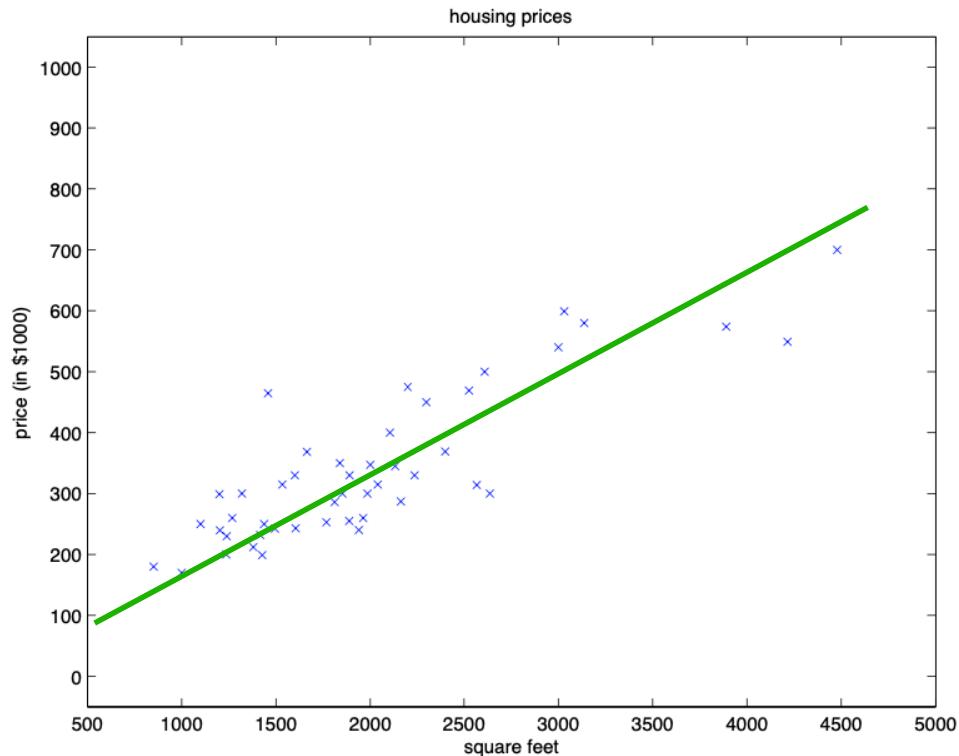Once we learned $\hat{w}$, we can use it to make prediction on any new feature $x$



housing prices

Given $x_{test}$, our prediction is:

$$\hat{y} = x_{test}^\top \hat{w}$$

$$= x_{test}^\top (XX^\top)^{-1}XY$$

# Prediction using linear regression

Once we learned $\hat{w}$, we can use it to make prediction on any new feature $x$



Given $x_{test}$, our prediction is:

$$\hat{y} = x_{test}^\top \hat{w}$$

$$= x_{test}^\top (XX^\top)^{-1} XY$$

$$= \sum_i \left( x_{test}^\top (XX^\top)^{-1} x_i \right) \cdot y_i$$

$n$

scalar

$= \sum_{i=1}^{n} \alpha_i y_i$

# Outline for Today

1. Intro on Linear Regression

2. Normal equation for linear Regression

3. Interpretation of Linear Regression using MLE / MAP

# Derive Linear regression via Maximum Likelihood Estimation

Assume $P(y \mid x; w) = \dfrac{1}{Z} \exp\left(-\dfrac{1}{2}(y - x^\top w)^2 / \sigma^2\right)$, i.e., $y = w^\top x + \epsilon,\ \epsilon \sim \mathcal{N}(0, \sigma^2)$

mean   $\sigma^2$

$$\prod_{i=1}^{n} P(y_i \mid x_i; w)$$

# Derive Linear regression via Maximum Likelihood Estimation

Assume $P(y \mid x; w) = \dfrac{1}{Z} \exp\left( -\dfrac{1}{2}(y - x^\top w)^2/\sigma^2 \right)$, i.e., $y = w^\top x + \epsilon, \ \epsilon \sim \mathcal{N}(0, \sigma^2)$

Let's maximize the log-likelihood of the data, i.e.,

# Derive Linear regression via Maximum Likelihood Estimation

Assume $P(y \mid x; w) = \dfrac{1}{Z} \exp\left(-\dfrac{1}{2}(y - x^\top w)^2 / \sigma^2\right)$, i.e., $y = w^\top x + \epsilon,\ \epsilon \sim \mathcal{N}(0, \sigma^2)$

Let's maximize the log-likelihood of the data, i.e.,

$$\arg\max_{w} \sum_{i=1}^{n} \ln P(y_i \mid x_i; w)$$

$$\ln\left(\frac{1}{Z} \exp\left(-\frac{1}{2}\left(y_i - x_i^\top w\right)^2 / \sigma^2\right)\right)$$

$$= \left(-\ln(Z)\right) - \frac{1}{2}\left(y_i - x_i^\top w\right)^\top / \sigma^2$$

# Derive Linear regression via Maximum Likelihood Estimation

Assume $P(y \mid x; w) = \dfrac{1}{Z} \exp\left(-\dfrac{1}{2}(y - x^\top w)^2 / \sigma^2\right)$, i.e., $y = w^\top x + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$

Let's maximize the log-likelihood of the data, i.e.,

$$\arg\max_{w} \sum_{i=1}^{n} \ln P(y_i \mid x_i; w)$$

$$= \arg\max_{w} \sum_{i=1}^{n} -\frac{1}{2\sigma^2}(w^\top x_i - y_i)^2 - \ln(Z)$$

$$\ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right)$$

$$Z := \sigma\sqrt{2\pi}$$

# Derive Linear regression via Maximum Likelihood Estimation

Assume $P(y \mid x; w) = \dfrac{1}{Z} \exp \left( -\dfrac{1}{2}(y - x^\top w)^2 / \sigma^2 \right)$, i.e., $y = w^\top x + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$
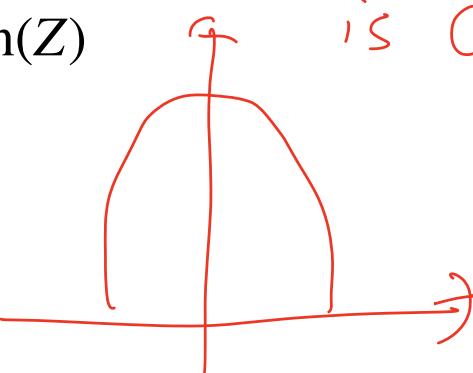
Let's maximize the log-likelihood of the data, i.e.,

$$\arg\max_{w} \sum_{i=1}^{n} \ln P(y_i \mid x_i; w)$$

$$= \arg\max_{w} \sum_{i=1}^{n} \cancel{\#} \frac{1}{2\sigma^2} (w^\top x_i - y_i)^2 - \ln(Z)$$

$$= \arg\min_{w} \sum_{i=1}^{n} (w^\top x_i - y_i)^2$$

$\ln P(\not\!y \mid x, w)$

is concave wrt $w$

# Derive Linear regression via MAP

Assume $P(y \mid x; w) = \dfrac{1}{Z} \exp\left(-\dfrac{1}{2}(y - x^\top w)^2 / \sigma^2\right)$, i.e., $y = w^\top x + \epsilon,\ \epsilon \sim \mathcal{N}(0, \sigma^2)$

To use MAP, we need to define a prior over $w$, we use Gaussian as well here:

$$w \sim \mathcal{N}(0, r^2 I)$$

$r$ big number

# Derive Linear regression via MAP

$$w \sim \mathcal{N}(0, r^2 I) \quad P(y \mid x; w) = \frac{1}{Z} \exp\left(-\frac{1}{2}(y - x^\top w)^2 / \sigma^2\right)$$

MAP:

# Derive Linear regression via MAP

$$w \sim \mathcal{N}(0, r^2 I) \qquad P(y \mid x; w) = \frac{1}{Z} \exp\left( -\frac{1}{2}(y - x^\top w)^2 / \sigma^2 \right)$$

MAP:

$$\arg\max_{w} \ln P(w \mid \mathcal{D}) \propto \text{prior} \times \text{likelihood}$$

# Derive Linear regression via MAP

$$w \sim \mathcal{N}(0, r^2 I) \qquad P(y \mid x; w) = \frac{1}{Z} \exp\left(-\frac{1}{2}(y - x^\top w)^2 / \sigma^2\right)$$

MAP:

$$\arg\max_{w} \ln P(w \mid \mathcal{D})$$

$$= \arg\max_{w} \ln P(w) + \ln P(\mathcal{D} \mid w)$$

# Derive Linear regression via MAP

$$w \sim \mathcal{N}(0, r^2 I) \qquad P(y \mid x; w) = \frac{1}{Z} \exp\left(-\frac{1}{2}(y - x^\top w)^2 / \sigma^2\right)$$

MAP:

$$\arg\max_{w} \ln P(w \mid \mathcal{D})$$

$$= \arg\max_{w} \ln P(w) + \ln P(\mathcal{D} \mid w)$$

$$= \arg\max_{w} \frac{-w^\top w}{2r^2} + \sum_{i=1}^{n} -\frac{1}{2\sigma^2}(w^\top x_i - y_i)^2$$

$$\hookrightarrow \ln P(w) \qquad \qquad \ln \prod_{i=1}^{n} P(y_i \mid x; w)$$

# Derive Linear regression via MAP

$$w \sim \mathcal{N}(0, r^2 I) \quad P(y \mid x; w) = \frac{1}{Z} \exp\left(-\frac{1}{2}(y - x^\top w)^2 / \sigma^2\right)$$

MAP:

$$\arg\max_{w} \ln P(w \mid \mathcal{D})$$

$$= \arg\max_{w} \ln P(w) + \ln P(\mathcal{D} \mid w)$$

$$= \arg\max_{w} \frac{-w^\top w}{2r^2} + \sum_{i=1}^{n} -\frac{1}{2\sigma^2}(w^\top x_i - y_i)^2$$

$$= \arg\min_{w} \frac{\sigma^2}{r^2} w^\top w + \sum_{i=1}^{n} (w^\top x_i - y_i)^2$$

# Derive Linear regression via MAP

$$w \sim \mathcal{N}(0, r^2 I) \qquad P(y \,|\, x; w) = \frac{1}{Z} \exp\left(-\frac{1}{2}(y - x^\top w)^2 / \sigma^2\right)$$
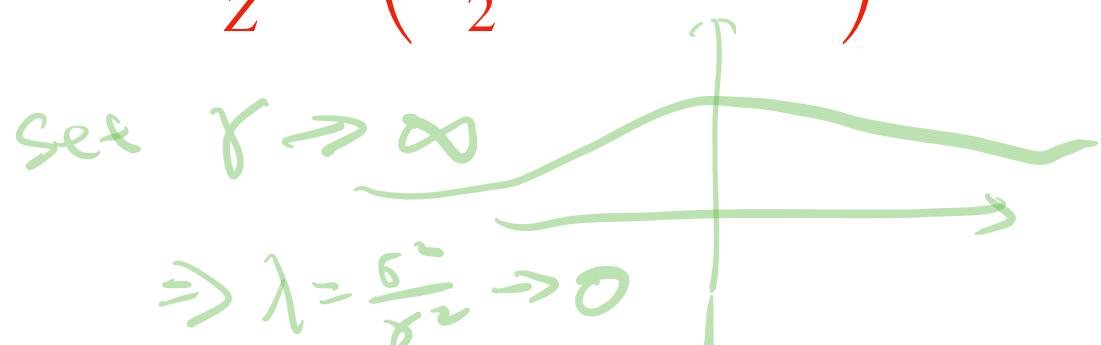
MAP:

$$\arg\max_{w} \ln P(w \,|\, \mathcal{D})$$

$$= \arg\max_{w} \ln P(w) + \ln P(\mathcal{D} \,|\, w)$$

$$= \arg\max_{w} \frac{-w^\top w}{2r^2} + \sum_{i=1}^{n} -\frac{1}{2\sigma^2}(w^\top x_i - y_i)^2$$

$$= \arg\min_{w} \frac{\sigma^2}{r^2} w^\top w + \sum_{i=1}^{n} (w^\top x_i - y_i)^2 \qquad = \arg\min_{w} \lambda \|w\|_2^2 + \sum_{i=1}^{n} (w^\top x_i - y_i)^2$$

set $\gamma \to \infty$

$$\Rightarrow \lambda = \frac{\sigma^2}{\gamma^2} \to 0$$

$$\lambda := \frac{\sigma^2}{\gamma^2}$$

Regularization

# Ridge Linear Regression

$$\arg \min_{w} \lambda \|w\|_2^2 + \sum_{i=1}^{n} (w^\top x_i - y_i)^2$$

In this case, we can derive a closed-form solution as well:

$$\hat{w} = (XX^\top + \lambda I)^{-1} XY$$

( Recall for normal Eqn:

$(XX^\top)^{-1} XY$

# Ridge Linear Regression

$$\arg \min_{w} \lambda \|w\|_2^2 + \sum_{i=1}^{n} (w^\top x_i - y_i)^2$$

In this case, we can derive a closed-form solution as well:

$$\hat{w} = (XX^\top + \lambda I)^{-1} XY$$

Note that it works even $XX^\top$ is not full rank

$XX^\top + \lambda I$

is always PD

$(\lambda > 0)$

# Summary for today

1. Linear regression, Normal equation, and MLE / MAP interpretation

# Summary for today

1. Linear regression, Normal equation, and MLE / MAP interpretation

2. Your take-home question: what is the SGD update rule for Linear regression? Is the update rule intuitively explainable?

# Summary for today

1. Linear regression, Normal equation, and MLE / MAP interpretation

2. Your take-home question: what is the SGD update rule for Linear regression? Is the update rule intuitively explainable?

3. Next Tue: Support Vector Machine!

$$f(x) = X^T A X \qquad A \in R^{d \times d}$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$\nabla_x f(x) = A X + A^T X$$

$$\nabla_x f(x) = \begin{bmatrix} \dfrac{\partial f(x)}{\partial x_1} \\[1em] \dfrac{\partial f(x)}{\partial x_2} \\[1em] \vdots \\[1em] \dfrac{\partial f(x)}{\partial x_d} \end{bmatrix} = A X + A^T X$$

$$A X = \sum_{i=1}^{d} A_i X_i$$

$$A = \begin{bmatrix} a_1 \cdots a_d \\ | \quad | \end{bmatrix}$$

$$X^T A X = X^T (A X)$$

$$= X^T \left( \sum_{i=1}^{d} a_i x_i \right)$$

$$= \sum_{i=1}^{d} \sum_{j=1}^{d} x_j a_{ij} x_i$$

$$\begin{bmatrix} | & | \\ a_1 & a_2 \\ | & | \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = a_1 x_1 + a_2 \cdot x_2$$

$$\underline{x^T A x} = \sum_{i=1}^{d} \sum_{j=1}^{d} x_j A_{i,j} x_i$$

$A \in \mathbb{R}^{d \times d}$

$$i \times \begin{bmatrix} - & A_{i,j} \\ & & \downarrow \\ & & \downarrow j \end{bmatrix} \qquad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$\frac{2 x^T A x}{2 x_1}$$