

The background of the slide features a large, faint watermark of the University of Turin seal. The seal is circular, with the Latin text "UNIVERSITAS STUDII TURINENSIS" around the perimeter and "MCCXXII" at the bottom. In the center is a shield depicting two figures, one seated and one standing, under a gothic arch.

Deep Learning

LM Computer Science, Data Science, Cybersecurity
2nd semester - 6 CFU

Luca Pasa, Nicolò Navarin & Alessandro Sperduti

Probability / Information theory

A primer on probability and information theory (chapter 3)

Maximum Likelihood estimation (section 5.5)

Probability

- **Random variable:** a variable that can take different values randomly
- Example: Tossing a coin: we could get Heads or Tails.
 - Heads=0 and Tails=1
 - In each experiment, Random Variable x can be either 0 or 1

*Random
Variable*

*Possible
Values*

*Random
Events*

$X =$

$\left\{ \begin{array}{l} 0 \\ 1 \end{array} \right.$

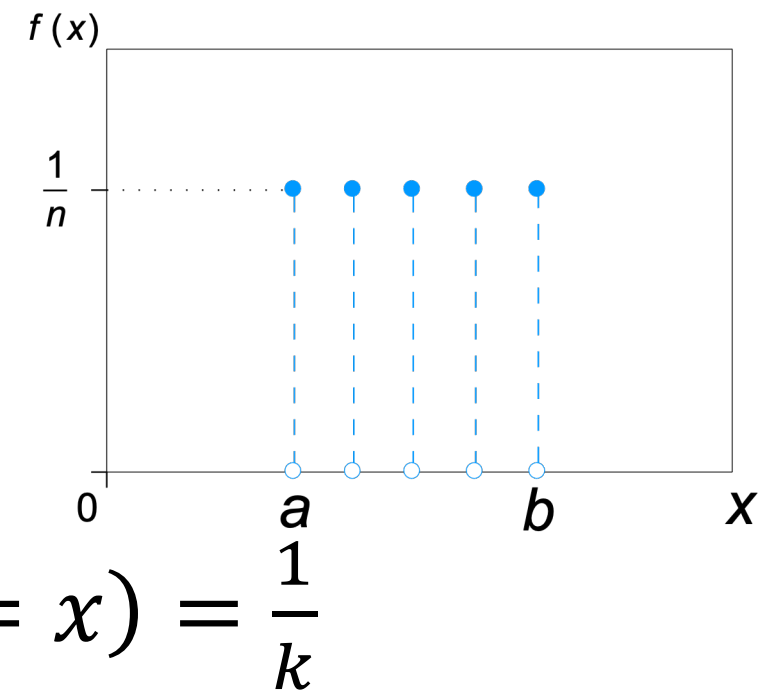
0

1



Probability Distributions

- **Probability Distribution**: A description of how likely a random variable x (or a set of random variables) is to take each of its possible states
- **Probability Mass Function** (Discrete variables)
 - Domain of P is the set of all possible states of x (k different values)
 - $\forall x \in \mathcal{X} \ 0 \leq P(x = x) \leq 1$
 - $\sum_{x \in \mathcal{X}} P(x) = 1$
- E.g. Uniform distribution $\forall_{x \in \mathcal{X}} P(x = x) = \frac{1}{k}$



Probability Distributions

- **Joint probability distribution**: probability distribution over 2 or more variables $P(x = x, y = y)$ or $P(x, y)$
 - Example 2 coins: $P(x = head, y = tail)$ or $P(head, tail)$
- **Marginalization**:

$$\forall x \in \mathbf{x} \ P(x = x) = \sum_y P(x = x, y = y)$$

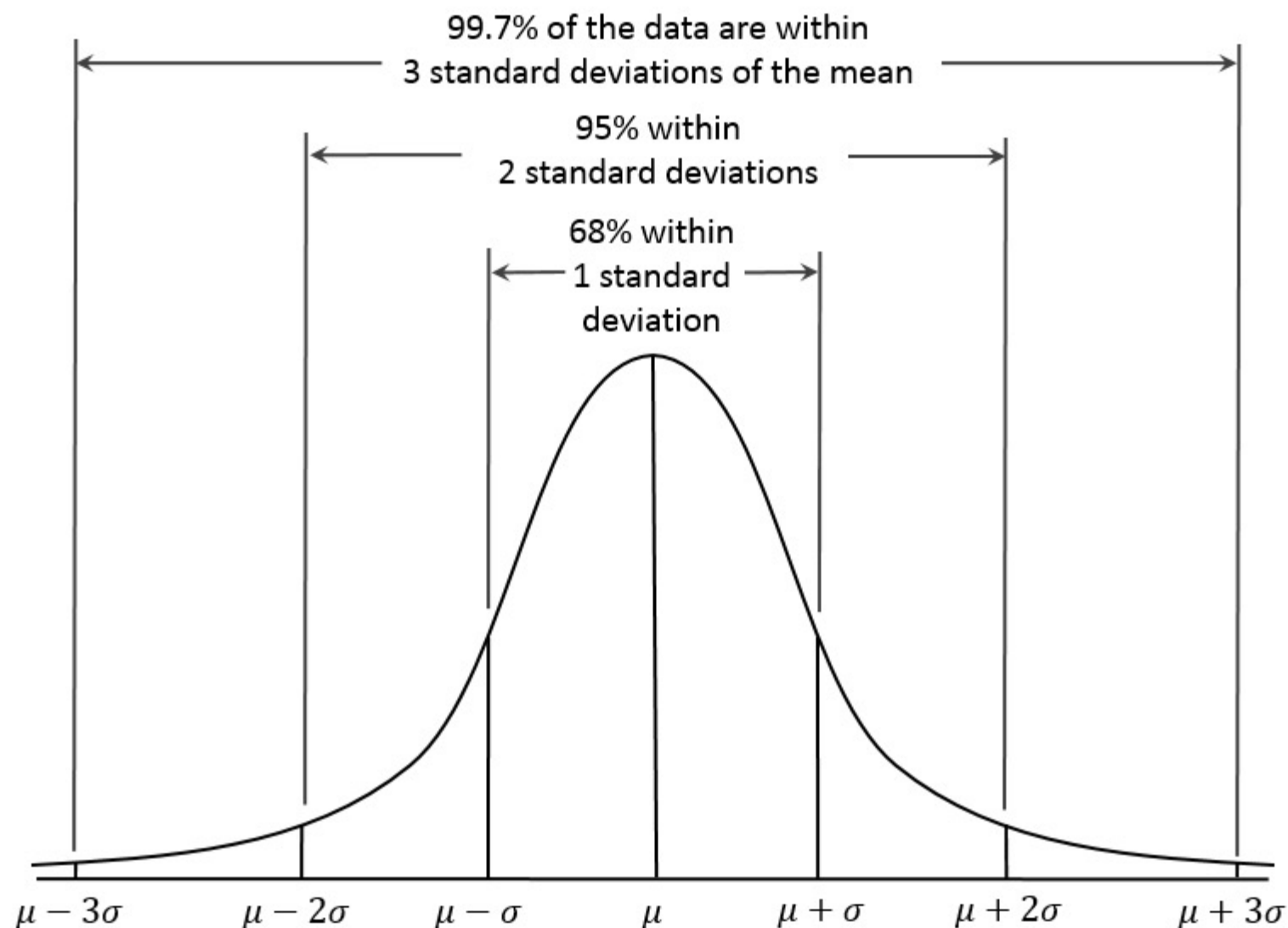
- **Continuous variables**: Probability distribution is described by a **Probability Density Function (PDF)**
 - Domain of p is the set of all possible states of x
 - $\forall x \in \mathbf{x} \ P(x) \geq 0$
 - $\int p(x) dx = 1$
- E.g. Gaussian distribution

Gaussian distribution

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

Parametrized by

PDF of
Gaussian
distribution



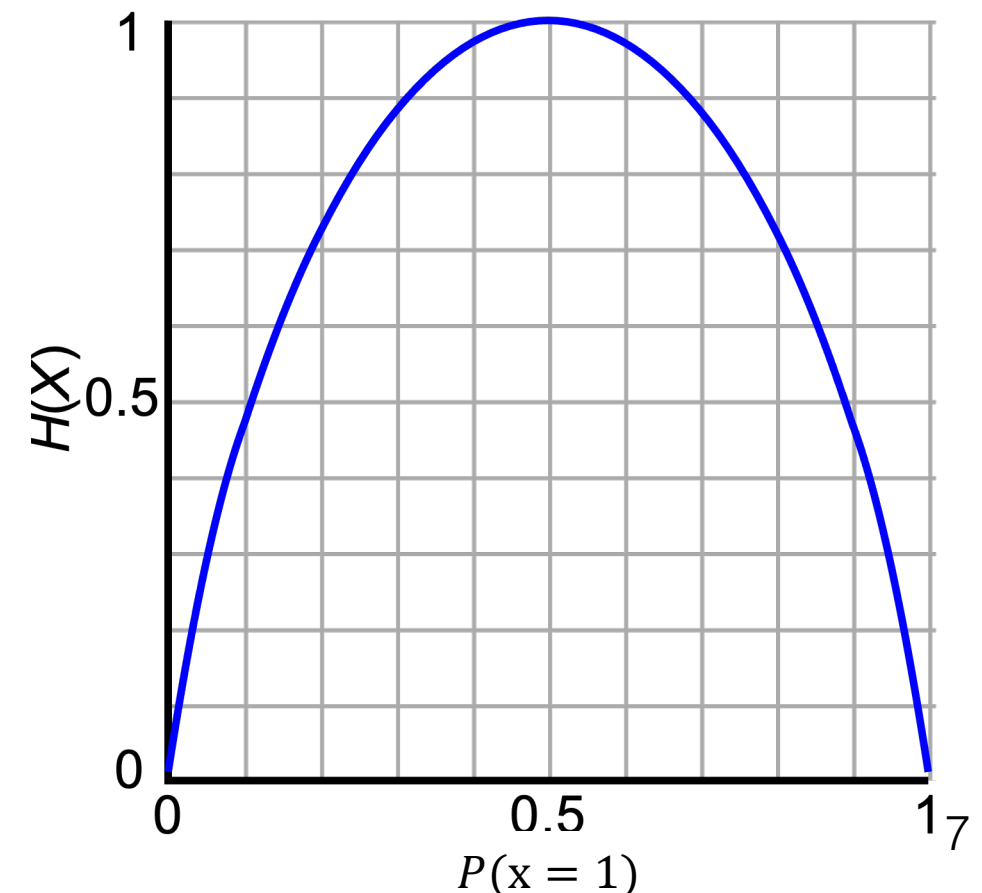
Entropy

- Shannon Entropy (discrete variable)

$$H(x) = -\mathbb{E}_{x \sim P(x)} [\log P(x)]$$

- Expected amount of (self-)information in an event drawn from distribution P
- Lower bound on the number of bits needed on average to encode a symbol drawn from that distribution

Entropy $H(X)$ of a coin flip, measured in bits, graphed versus the bias of the coin $P(x = 1)$, where $x = 1$ represents a result of heads.



Kullback-Leibler divergence and Cross Entropy

- Let's consider two probability distributions $P(x)$ and $Q(x)$

- Measure how different they are

$$D_{KL}(P \parallel Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]$$

- It is not a true distance because it is **not symmetric**


- $D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$


- Cross Entropy

$$H(P, Q) = H(P) + D_{KL}(P \parallel Q) = - \mathbb{E}_{x \sim P} [\log Q(x)]$$

Note: minimizing CE of P w.r.t. Q is equivalent to minimize KL divergence between P and Q (if P is given, $H(P)$ and $\mathbb{E}_{x \sim P} [\log P(x)]$ are constants)

Maximum likelihood estimation

- Principled way to derive estimators (models)
- Consider n examples $Tr = \{\mathbf{x}^1, \dots, \mathbf{x}^n\}$ drawn **i.i.d.** from $p_{data}(\mathbf{x})$  **Not known in advance**
- Let us consider a family of parametric probability distributions (models) $p_{model}(\mathbf{x}; \boldsymbol{\theta})$.
 - $p_{model}(\mathbf{x}; \boldsymbol{\theta})$ maps a point \mathbf{x} to a real number, **estimating** $p_{data}(\mathbf{x})$
 - Maximum Likelihood estimation for $\boldsymbol{\theta}$ is

$$\boldsymbol{\theta}_{ML} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p_{model}(Tr; \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i=1}^n p_{model}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$$


Here we do not know $p_{data}(\mathbf{x})$, we just know that $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$ are i.i.d!

Assumption: independence!
 $P(x, y) = P(x)P(y)$
If x and y are independent

..A side note on maximum likelihood

- ML is a special case of **maximum a posteriori estimation (MAP)**
- ML assumes a uniform prior distribution of the hypothesis
- MAP and maximum likelihood approach make predictions using a single point estimate of θ
- the Bayesian approach is to make predictions using a full probability distribution over θ

..A side note on maximum likelihood

Given a new instance x , what is the most probable *classification* ?

- ▶ $h_{MAP}(x)$, in general, is not the most probable classification!

Example: let's consider:

- ▶ three possible hypotheses:

$$P(h_1|D) = .4, \quad P(h_2|D) = .3, \quad P(h_3|D) = .3$$

- ▶ given a new instance x ,

$$h_1(x) = +, \quad h_2(x) = -, \quad h_3(x) = -$$

- ▶ what is the most probable classification for x ?

Maximum likelihood estimation

$$\boldsymbol{\theta}_{ML} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i=1}^n p_{model}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

- Taking the product of many probabilities is numerically unstable.
 - We can apply the log and the $\underset{\boldsymbol{\theta}}{\operatorname{argmax}}$ does not change (log-likelihood)

$$\begin{aligned}\boldsymbol{\theta}_{ML} &= \log(\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i=1}^n p_{model}(\mathbf{x}^{(i)}; \boldsymbol{\theta})) \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i=1}^n \log p_{model}(\mathbf{x}^{(i)}; \boldsymbol{\theta})\end{aligned}$$

Maximum likelihood estimation

$$\boldsymbol{\theta}_{ML} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i=1}^n \log p_{model}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

- We can equivalently divide by n to express ML as an expectation over training data

$$\boldsymbol{\theta}_{ML} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{data}} [\log p_{model}(\mathbf{x}; \boldsymbol{\theta})]$$

- ML minimizes the dissimilarity between \hat{p}_{data} and p_{model} , measured by the KL divergence

ML estimation as KL divergence

- KL divergence between \hat{p}_{data} and p_{model}

$$D_{KL}(\hat{p}_{data} \parallel p_{model}) = \mathbb{E}_{x \sim \hat{p}_{data}} [\log \hat{p}_{data}(x)] - \log p_{model}(x; \theta)$$

does not depend on the model

- To minimize the KL, we need only to minimize

$$\arg \min_{\theta} -\mathbb{E}_{x \sim \hat{p}_{data}} [\log p_{model}(x; \theta)]$$

- That is the same equation of ML in previous slide
- It also corresponds to minimizing the **cross-entropy** between the two distributions

Conditional Probability

- Probability of an event, given that some other event has happened.

The diagram illustrates Bayes' rule with the following components:

- LIKELIHOOD** (orange text): the probability of "B" being TRUE given that "A" is TRUE. An arrow points from this text to the $P(B|A)$ term in the numerator.
- PRIOR** (teal text): the probability of "A" being TRUE. An arrow points from this text to the $P(A)$ term in the numerator.
- POSTERIOR** (green text): the probability of "A" being TRUE given that "B" is TRUE. An arrow points from this text to the $P(A|B)$ term in the denominator.
- The probability of "B" being TRUE** (pink text): An arrow points from this text to the $P(B)$ term in the denominator.

The equation is written as:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Bayes' rule

Conditional log likelihood

- We can use ML to estimate a **conditional** probability $P(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$ to predict \mathbf{y} given \mathbf{x} (supervised learning)

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{Y} | \mathbf{X}; \boldsymbol{\theta})$$

- If input examples are i.i.d.

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \log P(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta})$$

Warning!

- We will use logarithm properties. Check Algebra cheat sheet if some of the rules applied in the next slide are not clear!

Algebra Cheat Sheet

Basic Properties & Facts

Arithmetic Operations

$$ab + ac = a(b + c)$$

$$\left(\frac{a}{b}\right) \div \frac{c}{d} = \frac{a}{bc}$$

$$\frac{a}{b} + \frac{c}{d} = \frac{ad + bc}{bd}$$

$$\frac{a-b}{c-d} = \frac{b-a}{d-c}$$

$$\frac{ab+ac}{a} = b+c, \quad a \neq 0$$

Exponent Properties

$$a^n a^m = a^{n+m}$$

$$(a^n)^m = a^{nm}$$

$$(ab)^n = a^n b^n$$

$$a^{-n} = \frac{1}{a^n}$$

$$\left(\frac{a}{b}\right)^{-n} = \left(\frac{b}{a}\right)^n = \frac{b^n}{a^n}$$

Properties of Radicals

$$\sqrt[n]{a} = a^{\frac{1}{n}} \quad \sqrt[n]{ab} = \sqrt[n]{a} \sqrt[n]{b}$$

$$\sqrt[m]{\sqrt[n]{a}} = \sqrt[mn]{a} \quad \sqrt[n]{\frac{a}{b}} = \frac{\sqrt[n]{a}}{\sqrt[n]{b}}$$

$$\sqrt[n]{a^n} = a, \text{ if } n \text{ is odd}$$

$$\sqrt[n]{a^n} = |a|, \text{ if } n \text{ is even}$$

Properties of Inequalities

If $a < b$ then $a + c < b + c$ and $a - c < b - c$

If $a < b$ and $c > 0$ then $ac < bc$ and $\frac{a}{c} < \frac{b}{c}$

If $a < b$ and $c < 0$ then $ac > bc$ and $\frac{a}{c} > \frac{b}{c}$

Properties of Absolute Value

$$|a| = \begin{cases} a & \text{if } a \geq 0 \\ -a & \text{if } a < 0 \end{cases}$$

$$|a| \geq 0 \quad |-a| = |a|$$

$$|ab| = |a||b| \quad \left|\frac{a}{b}\right| = \frac{|a|}{|b|}$$

$$|a+b| \leq |a| + |b| \quad \text{Triangle Inequality}$$

Distance Formula

If $P_1 = (x_1, y_1)$ and $P_2 = (x_2, y_2)$ are two points the distance between them is

$$d(P_1, P_2) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Complex Numbers

$$i = \sqrt{-1} \quad i^2 = -1 \quad \sqrt{-a} = i\sqrt{a}, \quad a \geq 0$$

$$(a+bi) + (c+di) = a+c + (b+d)i$$

$$(a+bi) - (c+di) = a-c + (b-d)i$$

$$(a+bi)(c+di) = ac - bd + (ad+bc)i$$

$$(a+bi)(a-bi) = a^2 + b^2$$

$$|a+bi| = \sqrt{a^2 + b^2} \quad \text{Complex Modulus}$$

$$\overline{(a+bi)} = a-bi \quad \text{Complex Conjugate}$$

$$(a+bi)(a+bi) = |a+bi|^2$$

Logarithms and Log Properties

Definition

$y = \log_b x$ is equivalent to $x = b^y$

Example

$$\log_5 125 = 3 \quad \text{because } 5^3 = 125$$

Special Logarithms

$\ln x = \log_e x$ natural log

$\log x = \log_{10} x$ common log

where $e = 2.718281828K$

Logarithm Properties

$$\log_b b = 1 \quad \log_b 1 = 0$$

$$\log_b b^x = x \quad b^{\log_b x} = x$$

$$\log_b (x^r) = r \log_b x$$

$$\log_b (xy) = \log_b x + \log_b y$$

$$\log_b \left(\frac{x}{y}\right) = \log_b x - \log_b y$$

The domain of $\log_b x$ is $x > 0$

Factoring and Solving

Factoring Formulas

$$x^2 - a^2 = (x+a)(x-a)$$

$$x^2 + 2ax + a^2 = (x+a)^2$$

$$x^2 - 2ax + a^2 = (x-a)^2$$

$$x^2 + (a+b)x + ab = (x+a)(x+b)$$

$$x^3 + 3ax^2 + 3a^2x + a^3 = (x+a)^3$$

$$x^3 - 3ax^2 + 3a^2x - a^3 = (x-a)^3$$

$$x^3 + a^3 = (x+a)(x^2 - ax + a^2)$$

$$x^3 - a^3 = (x-a)(x^2 + ax + a^2)$$

$$x^{2n} - a^{2n} = (x^n - a^n)(x^n + a^n)$$

If n is odd then,

$$x^n - a^n = (x-a)(x^{n-1} + ax^{n-2} + \dots + a^{n-1})$$

$$x^n + a^n = (x+a)(x^{n-1} - ax^{n-2} + a^2x^{n-3} - \dots + a^{n-1})$$

Completing the Square

$$\text{Solve } 2x^2 - 6x - 10 = 0$$

(1) Divide by the coefficient of the x^2

$$x^2 - 3x - 5 = 0$$

(2) Move the constant to the other side.

$$x^2 - 3x = 5$$

(3) Take half the coefficient of x , square it and add it to both sides

$$x^2 - 3x + \left(-\frac{3}{2}\right)^2 = 5 + \left(-\frac{3}{2}\right)^2 = 5 + \frac{9}{4} = \frac{29}{4}$$

(4) Factor the left side

$$\left(x - \frac{3}{2}\right)^2 = \frac{29}{4}$$

(5) Use Square Root Property

$$x - \frac{3}{2} = \pm \sqrt{\frac{29}{4}} = \pm \frac{\sqrt{29}}{2}$$

(6) Solve for x

$$x = \frac{3}{2} \pm \frac{\sqrt{29}}{2}$$

Quadratic Formula

Solve $ax^2 + bx + c = 0, \quad a \neq 0$

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

If $b^2 - 4ac > 0$ - Two real unequal solns.

If $b^2 - 4ac = 0$ - Repeated real solution.

If $b^2 - 4ac < 0$ - Two complex solutions.

Square Root Property

If $x^2 = p$ then $x = \pm \sqrt{p}$

Absolute Value Equations/Inequalities

If b is a positive number

$$|p| = b \Rightarrow p = -b \quad \text{or} \quad p = b$$

$$|p| < b \Rightarrow -b < p < b$$

$$|p| > b \Rightarrow p < -b \quad \text{or} \quad p > b$$

Linear Regression as Maximum Likelihood

- Linear Regression: algorithm that learns to take an input x and produce an output value \hat{y}
 - minimize **mean squared error**, why??
- Let's revisit linear regression from the point of view of maximum likelihood estimation
 - Define a linear regression model that produces a conditional distribution $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(y; \hat{y}(\mathbf{x}; \boldsymbol{\theta}), \sigma^2)$
 - output $\hat{y}(\mathbf{x}; \boldsymbol{\theta})$, the **mean of a Gaussian distribution**
 - (variance is fixed to some constant)

Linear Regression as Maximum Likelihood

- Examples are assumed to be i.i.d

$$\begin{aligned}\boldsymbol{\theta}_{ML} &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) = \\ &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \mathcal{N}(y^{(i)}; \hat{y}(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \sigma^2) = \\ &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \log \left(\sqrt{\frac{1}{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y^{(i)} - \hat{y}^{(i)})^2} \right)\end{aligned}$$

Maximise the
log likelihood

Gaussian
distribution

$\boldsymbol{\theta}$: parameters of
the linear
regression that
computes $\hat{y}^{(i)}$

- $\hat{y}^{(i)}$ is the output of the linear regression $\hat{y}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$

$$\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^n \log \left(\sqrt{\frac{1}{2\pi\sigma^2}} e^{\left(-\frac{1}{2\sigma^2}(y^{(i)} - \hat{y}^{(i)})^2\right)} \right) =$$

Log product rule

$$= \arg \max_{\theta} \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) + \log e^{\left(-\frac{1}{2\sigma^2}(y^{(i)} - \hat{y}^{(i)})^2\right)} =$$

Log quotient rule

$$= \arg \max_{\theta} \sum_{i=1}^n \log(1) - \log(\sqrt{2\pi\sigma^2}) + \log e^{\left(-\frac{1}{2\sigma^2}(y^{(i)} - \hat{y}^{(i)})^2\right)}$$

Log power rule

$$= \arg \max_{\theta} \sum_{i=1}^n \log(1) - \log(\sqrt{2\pi\sigma^2}) - \frac{1}{2\sigma^2} (y^{(i)} - \hat{y}^{(i)})^2 \log(e)$$

Log power rule

$$= \arg \max_{\theta} \sum_{i=1}^n -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y^{(i)} - \hat{y}^{(i)})^2$$

Algebra

$$= \arg \max_{\theta} -\frac{n}{2} \log(2\pi\sigma^2) + \sum_{i=1}^n -\frac{1}{2} \left(\frac{(y^{(i)} - \hat{y}^{(i)})^2}{\sigma^2} \right)$$

$$= \arg \max_{\theta} -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

Comparing ML with MSE

$$\theta_{ML} = \arg \max_{\theta} \underbrace{-\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2}_{\text{Does not depend on } \theta} =$$

$$= \arg \max_{\theta} - \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

$$\theta_{MSE} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

maximizing the log-likelihood with respect to θ yields the same estimate of the parameters θ as does minimizing the mean squared error.