**Computer Science**

# CSC380: Principles of Data Science
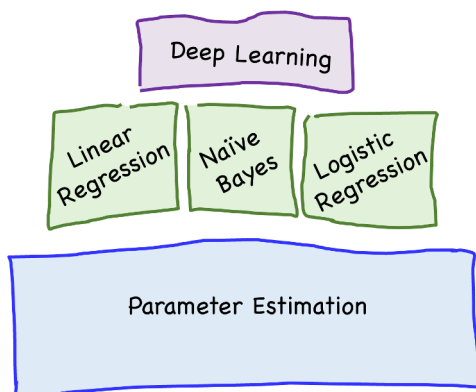
## Statistics 2

Credit:
- Jason Pacheco,
- Kwang-Sung Jun,
- Chicheng Zhang
- Xinchen yu

1

## Our path

2

Deep Learning

Linear Regression   Naïve Bayes   Logistic Regression

Parameter Estimation

- We don't know the true parameter.

- But we have observations.

- We assume each i.i.d observation follows a probability distribution with unknown parameters, and we build model.
  - e.g., Naive bayes model (X~Bernoulli)

- Compute estimator to estimate true parameter

- Many types of estimators with different properties
  - consistency
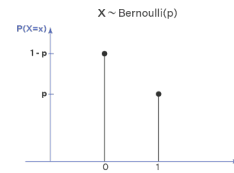  - efficiency (mean squared error)
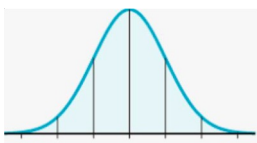  - unbiasedness

2

## Review

**Law of Large Numbers:**          [1, 0, 1, 0, 0, …, 1, 1, 0]          (1+0+1+0+0+...+1+1+0)/N

$$\lim_{N\to\infty} \hat{\mu}_n = \mu$$

Draw from a distribution with unknown mean

$X \sim \text{Bernoulli}(p)$

**Central Limit Theorem:**

[1, 0, 1, 0, 0, …, 1, 0, 1]  $\bar{X}_N$ for sample 1

[1, 0, 0, 0, 0, …, 1, 1, 0]  $\bar{X}_N$ for sample 2

[1, 1, 1, 0, 1, …, 0, 1, 0]  $\bar{X}_N$ for sample 3

……

$$\lim_{N\to\infty} \bar{X}_N \to \mathcal{N}\left(\mu, \frac{\sigma^2}{N}\right)$$

$$\lim_{N\to\infty} \frac{\sqrt{N}}{\sigma}(\bar{X}_N - \mu) \to \mathcal{N}(0,1)$$

[0, 0, 1, 1, 1, …, 0, 0, 0]  $\bar{X}_N$ for sample k

If N is very large, and we draw the distribution of $\bar{X}_N$ from all the samples, it follows normal distribution.

3

## Intuition Check

*Suppose that we toss a coin 100 times.  We observe 73 heads and 27 tails…*

<u>Question</u> Let $\theta$ be the coin bias (probability of heads).  What is a more likely estimate?  What is your reasoning?

**A:** $\hat{\theta} = 0.73$, strong preference for heads          Why sample mean?

**B:** $\hat{\theta} = 0.50$, fair coin (we observed unlucky outcomes)

**Likelihood (informally)** Probability/density of the observed outcomes from a particular model.
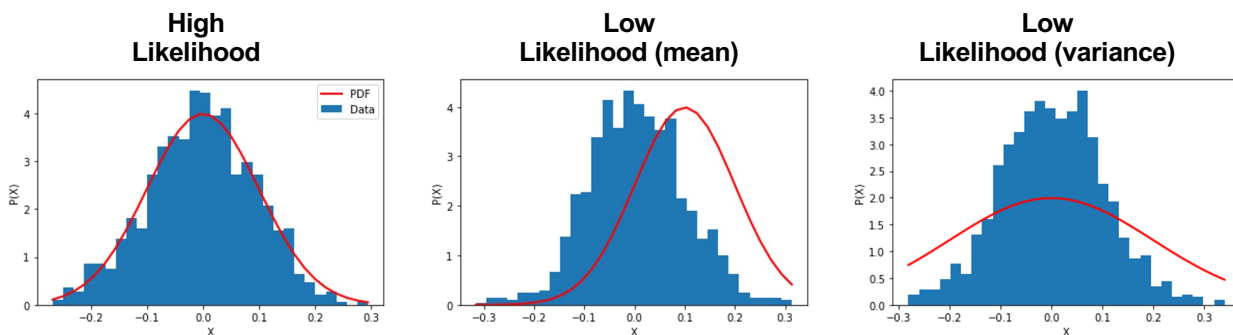
4

## Likelihood (Intuitively)

*Suppose we observe N data points from a Gaussian model $\mathcal{N}(\mu, \sigma^2)$, and wish to estimate both $\mu$ and $\sigma^2$.*

*Say we only need to choose from the following three Gaussians…*



***Likelihood Principle:*** *Given a statistical model, the <u>likelihood function </u>describes all evidence of a parameter that is contained in the data.*

5

## Likelihood Function

Suppose $x_i \sim p(x; \theta)$, then what is the **joint probability** over N *independent identically distributed* (iid) observations $x_1, \ldots, x_N$?

$$p(x_1, \ldots, x_N; \theta) = \prod_{i=1}^{N} p(x_i; \theta)$$

what appears after ; are parameters, not random variables.

- We call this the **likelihood function**, often denoted $\mathcal{L}_N(\theta)$
- It is a function of the parameter $\theta$, the data are fixed
- Describes how well parameter $\theta$ describes data (<u>*goodness of fit*</u>)

*How could we use this to estimate a parameter $\theta$?*

6

## Likelihood Function

Suppose $x_i \sim p(x; \theta)$, then what is the **joint probability** over N *independent identically distributed* (iid) observations $x_1, \ldots, x_N$?

$$p(x_1, \ldots, x_N; \theta) = \prod_{i=1}^{N} p(x_i; \theta)$$

what appears after ; are parameters, not random variables.

Suppose X ~ Bernoulli(p), we have 5 observations [1, 1, 0, 1, 0]

- If true parameter is **0.6**: fit the data better

$p(1,\ 1,\ 0,\ 1,\ 0;\ .6) = p(1; .6) \cdot p(1; .6) \cdot p(0; .6) \cdot p(1; .6) \cdot p(0; .6) = 0.6^3 0.4^2$  =.03

- If true parameter is 0.2:

$p(1,\ 1,\ 0,\ 1,\ 0;\ .2) = p(1; .2) \cdot p(1; .2) \cdot p(0; .2) \cdot p(1; .2) \cdot p(0; .2) = 0.2^3 0.8^2$ =.01

## Maximum Likelihood

**Maximum Likelihood Estimator (MLE)** as the name suggests, maximizes the likelihood function.

$$\hat{\theta}^{\mathrm{MLE}} = \arg\max_{\theta} \mathcal{L}_N(\theta) = \prod_{i=1}^{N} p(x_i; \theta)$$

**Question** How do we find the MLE?

1. closed-form
2. iterative methods
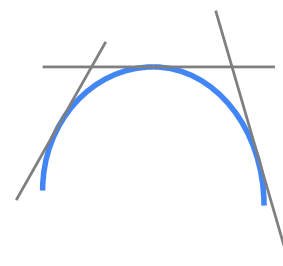
## How to find the maximum/maximizer of a function?

9

## Option 1: finding the maximum/maximizer

Example: Suppose $f(\theta) = -a\theta^2 + b\theta + c$ with $a > 0$

It is a quadratic function.
=> finding the 'flat' point suffices

Compute the gradient and set it equal to 0
$$f'(\theta) = -2a\theta + b \qquad => \quad \theta = \frac{b}{2a}$$
Closed form!

Q: Does this trick of gradient=0 work for other functions?
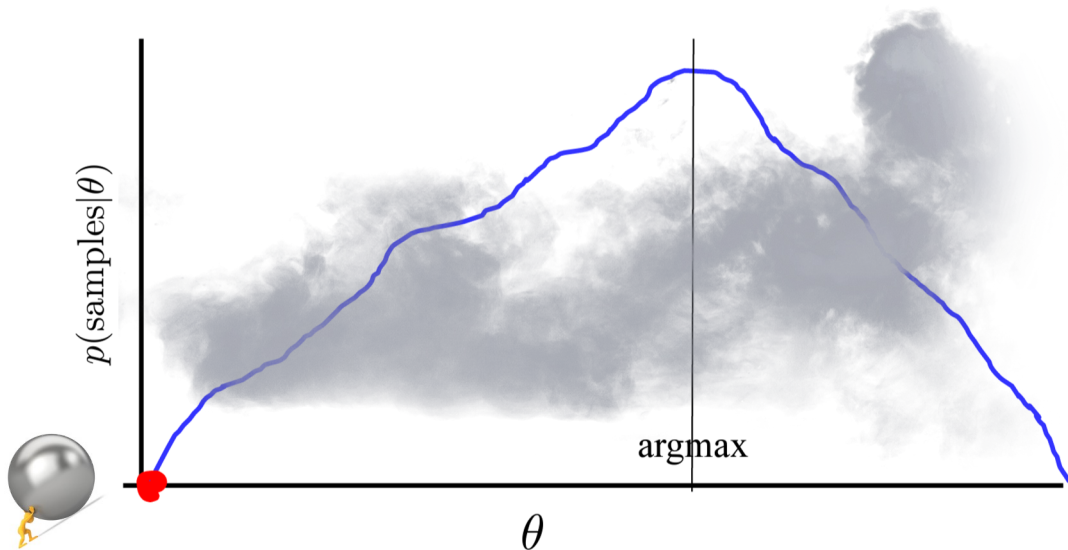⇒ Yes, **concave** functions!
⇒ This method finds extreme points – could be maximum and minimum.
⇒ For concave function, a local maximum is also a global one
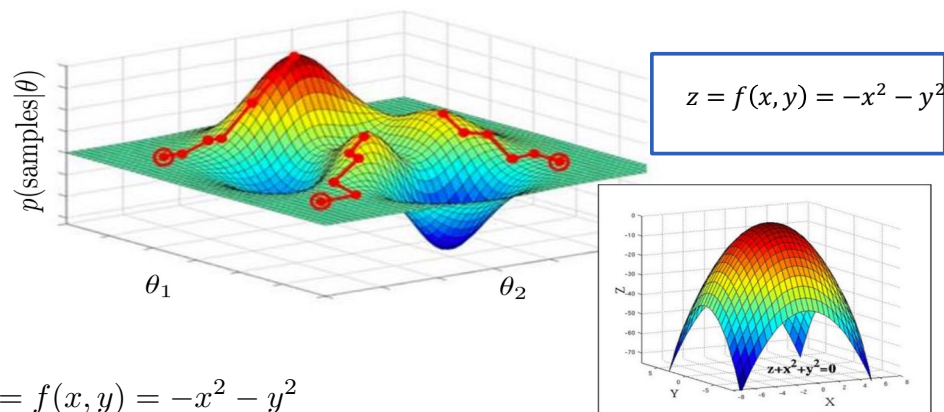
10

## Option 2: finding the maximum/maximizer

Walk uphill and you will find a local maxima (if your step size is small enough)

11

## Option 2: finding the maximum/maximizer

$$z = f(x,y) = -x^2 - y^2$$

Here $z = f(x,y) = -x^2 - y^2$

$$\frac{\partial}{\partial x} f(x,y) = -2x \qquad \text{and} \qquad \frac{\partial}{\partial x} f(x,y) = -2y$$

So starting say at $p_1 = (1,1)$, walk in direction $\vec{d} = (-2,-2)$
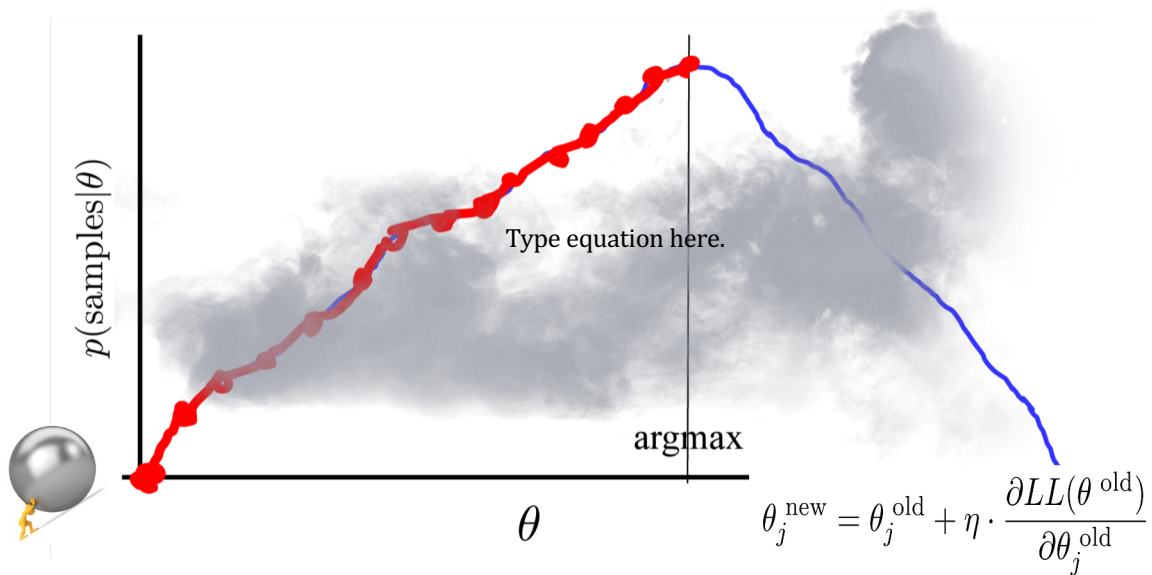next point. Say $\eta = 0.1$
$p_2 = (1,1) + \eta\vec{v} = (1,1) + \eta(-2,-2) = (1.9,1.9)$

Walk uphill and you will find a local maxima (if your step size is small enough)

12

## Option 2: finding the maximum/maximizer

Type equation here.

argmax

$\theta$

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \frac{\partial LL(\theta^{\text{old}})}{\partial \theta_j^{\text{old}}}$$
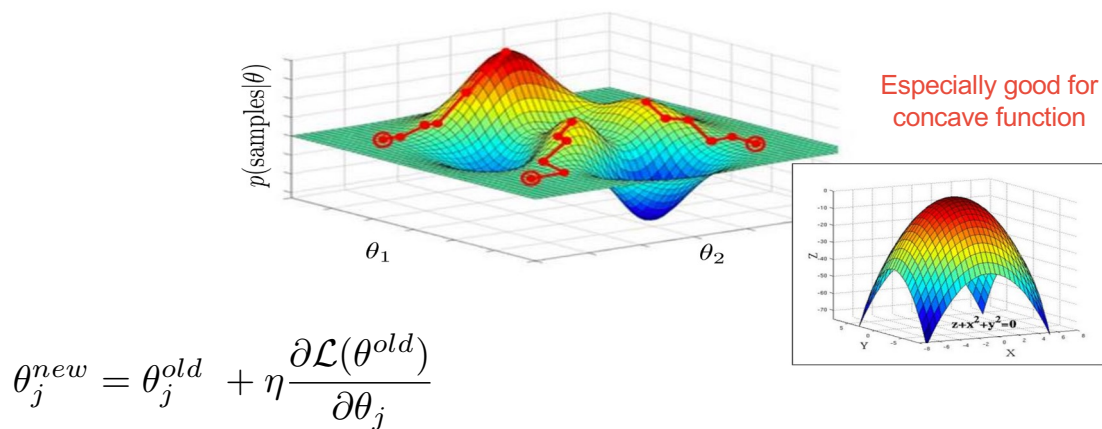
Walk uphill and you will find a local maxima (if your step size is small enough)
Here $\eta$ is the step size – if too large, we could overshoot.

13

## Option 2: finding the maximum/maximizer

Especially good for
concave function

$\theta_1$    $\theta_2$

$$\theta_j^{new} = \theta_j^{old} + \eta \frac{\partial \mathcal{L}(\theta^{old})}{\partial \theta_j}$$

Walk uphill and you will find a local maxima (if your step size is small enough)

14

## Option 2: finding the maximum/maximizer

What if there is no closed form solution?

Example: $f(x) = \frac{1}{2}x(ax - 2\log(x) + 2)$

$f'(x) = ax - \log(x)$

No known closed form for $ax = \log(x)$

Iterative methods:
- Gradient ascent (or *descent* if you are minimizing):
- Newton's method
- Etc. (beyond the scope of our class)

Iterative methods
- for **concave** functions
  => Will find the global maximum
- for **nonconcave**,
  => usually find a local maximum but could also get stuck at *stationary point*.

$$\theta_j^{new} = \theta_j^{old} + \eta \frac{\partial \mathcal{L}(\theta^{old})}{\partial \theta_j}$$
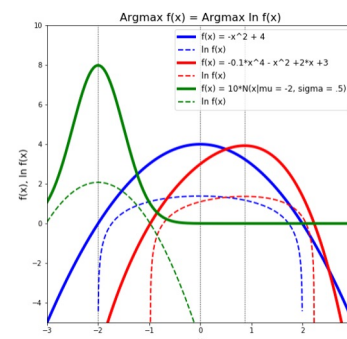
## Maximum Likelihood

Maximizing **log**-likelihood makes the math easier (as we will see) and doesn't change the answer (logarithm is an increasing function)

$$\hat{\theta}^{\mathrm{MLE}} = \arg\max_{\theta} \ \log \mathcal{L}_N(\theta) = \sum_{i=1}^{N} \log p(x_i; \theta)$$

Derivative is a linear operator so,

$$\frac{d}{d\theta} \log \mathcal{L}_N(\theta) = \sum_{i=1}^{N} \frac{d}{d\theta} \log p(x_i; \theta)$$

One term per data point
Can be computed in parallel
(big data)

## Review: maximum likelihood estimation

1. Decide on a model for the likelihood of your samples. This is often using a PMF or PDF.

2. Write out the log likelihood function.

3. State that the optimal parameters are the argmax of the log likelihood function.

4. Calculate the derivative of LL with respect to theta

5. Use an optimization algorithm to calculate argmax

17

## Maximum Likelihood: Bernoulli

**Example**: Consider I.I.D. random variables: $X_1$, $X_2$, $X_3$ ... $X_n \sim Bernoulli\,(p)$
We don't know the coin bias $p$.

Probability Mass function: $\quad p^{x_i}(1-p)^{1-x_i}$

Likelihood: $\displaystyle \mathcal{L}_n(p) = \prod_{i=1}^{n} p^{x_i}(1-p)^{1-x_i} = p^{x_1+\ldots+x_n}(1-p)^{n-(x_1+\ldots+x_n)}$ $\qquad \boxed{S = \sum_i x_i}$

$$= p^S(1-p)^{n-S}$$

Log likelihood: $\quad \mathcal{LL}_n(p) = S \log p + (n-S)\log(1-p)$

18

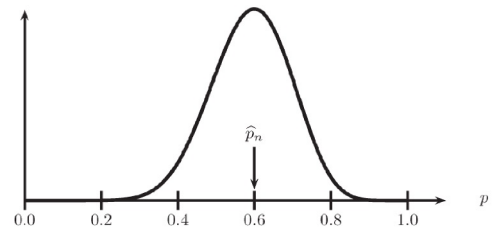## Maximum Likelihood: Bernoulli

Set the derivative of $\mathcal{LL}_n(p)$ to zero and solve,

$$\mathcal{LL}_n(p) = S \log p + (n - S)\log(1 - p)$$

$$\frac{\partial \mathcal{LL}_n(p)}{\partial p} = S\frac{1}{p} + (n - S)\frac{-1}{1 - p} = 0$$



*Likelihood function for Bernoulli with n=20 and $\sum_i x_i = 12$ heads*

We get:

$$p_{MLE} = \frac{S}{n} = \frac{1}{n}\sum_i x_i$$

$$\boxed{S = \sum_i x_i}$$

Isn't that the same as the sample mean?

Yes, for Bernoulli

$\Rightarrow$ this showcases how MLE is aligned to our intuition!

## Maximum Likelihood: Gaussian

**Example** Let $X_1, \ldots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ with parameters $\theta = (\mu, \sigma^2)$ and the likelihood function (ignoring some constants) is:

$$
\begin{aligned}
\mathcal{L}_n(\mu, \sigma) &= \prod_i \frac{1}{\sigma} \exp\left\{ -\frac{1}{2\sigma^2}(X_i - \mu)^2 \right\} \\
&= \sigma^{-n} \exp\left\{ -\frac{1}{2\sigma^2}\sum_i (X_i - \mu)^2 \right\} \\
&= \sigma^{-n} \exp\left\{ -\frac{nS^2}{2\sigma^2} \right\} \exp\left\{ -\frac{n(\overline{X} - \mu)^2}{2\sigma^2} \right\}
\end{aligned}
$$

Could show (next slide) – algebra

$$\sum_i (X_i - \mu)^2 = nS^2 + n(\bar{X} - \mu)^2$$

$$\boxed{e^{x+y} = e^x e^y}$$

Where $\bar{X} = \frac{1}{n}\sum_i X_i$ and $S^2 = \frac{1}{n}\sum_i (X_i - \bar{X})^2$ are sample mean and sample variance, respectively.

## Maximum Likelihood: Gaussian

$$\sum_i (X_i - \mu)^2 = \sum_i (X_i - \overline{X} + \overline{X} - \mu)^2 = \sum_i \left[ (X_i - \overline{X})^2 + 2(X_i - \overline{X})(\overline{X} - \mu) + (\overline{X} - \mu)^2 \right]$$

$$= \sum_i \left[ (X_i - \overline{X})^2 + 2\left( X_i\overline{X} - X_i\mu - \overline{X}^2 + \overline{X}\mu \right) + \left( \overline{X}^2 - 2\overline{X}\mu + \mu^2 \right) \right]$$

Given:

$\bar{X} = \frac{1}{n} \sum_i X_i$

$S^2 = \frac{1}{n} \sum_i (X_i - \bar{X})^2$

$$= \sum_i \left[ (X_i - \overline{X})^2 + 2X_i\overline{X} - 2X_i\mu - 2\overline{X}^2 + 2\overline{X}\mu + \overline{X}^2 - 2\overline{X}\mu + \mu^2 \right]$$

$$= \sum_i \left[ (X_i - \overline{X})^2 + 2X_i(\overline{X} - \mu) - \overline{X}^2 + \mu^2 \right]$$

$$= \sum_i (X_i - \overline{X})^2 + \sum_i 2X_i(\overline{X} - \mu) - n\overline{X}^2 + n\mu^2)$$

$$= \sum_i (X_i - \overline{X})^2 + 2n\overline{X}(\overline{X} - \mu) - n\overline{X}^2 + n\mu^2$$

$$= \sum_i (X_i - \overline{X})^2 + n\left( \overline{X}^2 - 2\overline{X}\mu + \mu^2 \right) = \sum_i (X_i - \overline{X})^2 + n(\overline{X} - \mu)^2$$

21

## Maximum Likelihood: Gaussian

$$\mathcal{L}(\mu, \sigma) = \sigma^{-n} \exp\left\{ -\frac{nS^2}{2\sigma^2} \right\} \exp\left\{ -\frac{n(\overline{X} - \mu)^2}{2\sigma^2} \right\}$$

Continuing, write log-likelihood as:

$$\ell(\mu, \sigma) = -n \log \sigma - \frac{nS^2}{2\sigma^2} - \frac{n(\overline{X} - \mu)^2}{2\sigma^2}.$$

Solve zero-gradient conditions:

$$\frac{\partial \ell(\mu, \sigma)}{\partial \mu} = 0 \quad \text{and} \quad \frac{\partial \ell(\mu, \sigma)}{\partial \sigma} = 0,$$

To obtain maximum likelihood estimates of mean / variance:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} X_i \qquad\qquad \hat{\sigma}^2 = \frac{1}{n} \sum_i (X_i - \hat{\mu})^2$$

22

## MLE Review: Probability/Density vs Likelihood

- The **probability/density** of data given parameter is mathematically the same object as **likelihood** of a parameter given data

- The difference is the <u>point of view</u>!
  - From the *probabilistic perspective*, the parameter is fixed and **PMF/PDF** is viewed as a function of the possible data

  - From the *statistical perspective*, the data is given (thus fixed) and we view **likelihood** as a function of the parameter.

- Statistics is inherently about reverse engineering.

23