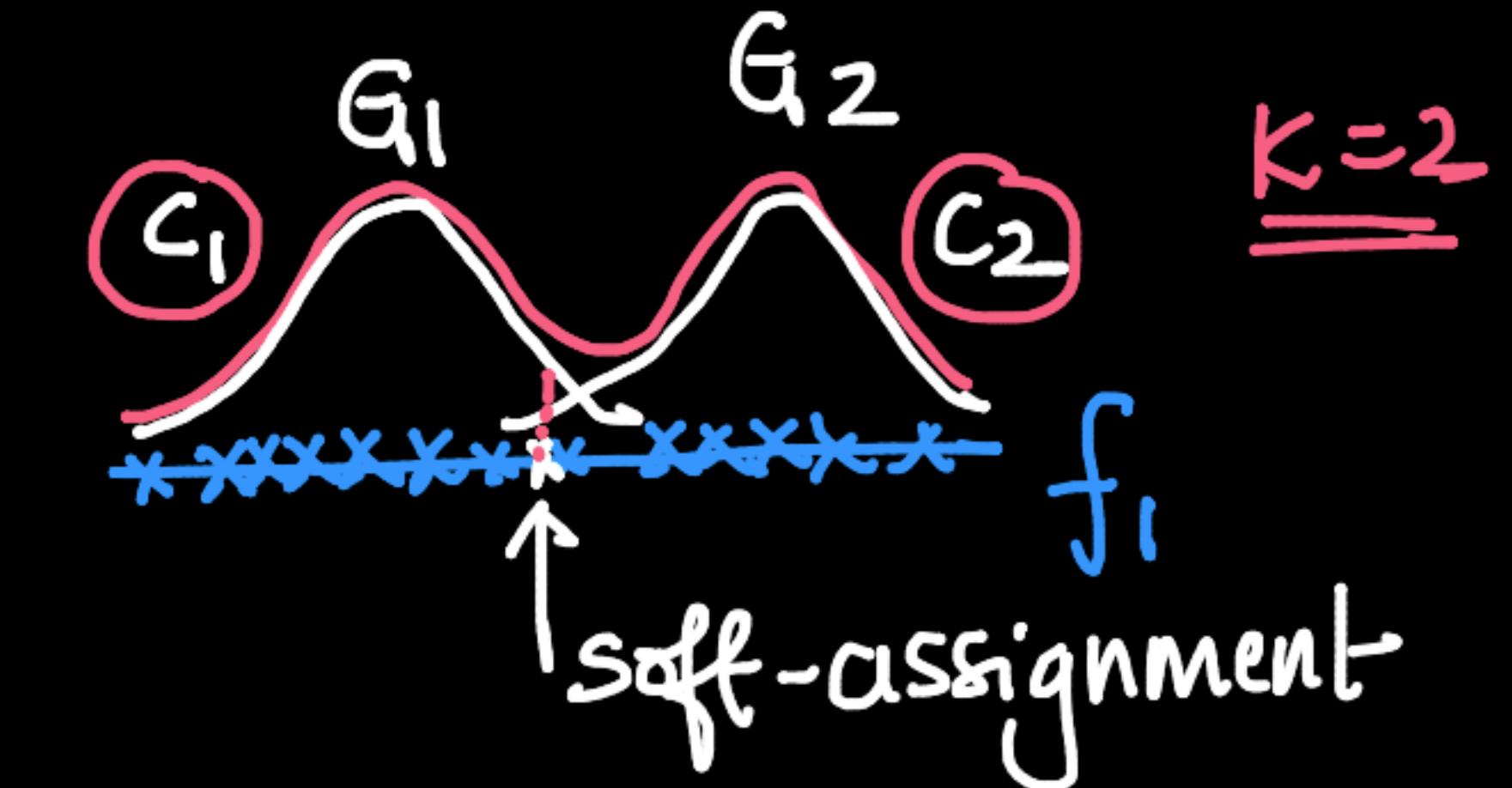


- ~~assumption~~ → Topics:
- GMM - optimization
 - =
- GMM: EM algorithm
- SUM → Kernel tricks
- { - Misc. topics - 1
 - ✓ - Lasso, Ridge & ElasticNet
 - ✓ - Calibration: Platt's scaling; Isotonic Regression
 - K-Means++ (Speedup)
 - Robust models: RANSAC
- ↑

Recap. of GMM: → Math-heavy

- ① Mixture-model →
- ② 1D-GMM
- ③ Multi-dimensional Gaussian: $N_d(\mu_d, \Sigma_{d \times d})$

1D Gaussian: $\underline{N(\mu, \sigma)}$



Multi-dimensional Gaussian: $N_d(\mu_d, \Sigma_{d \times d})$

↳ hill in $d+1$ -dim ← PDF

Cov-Matrix

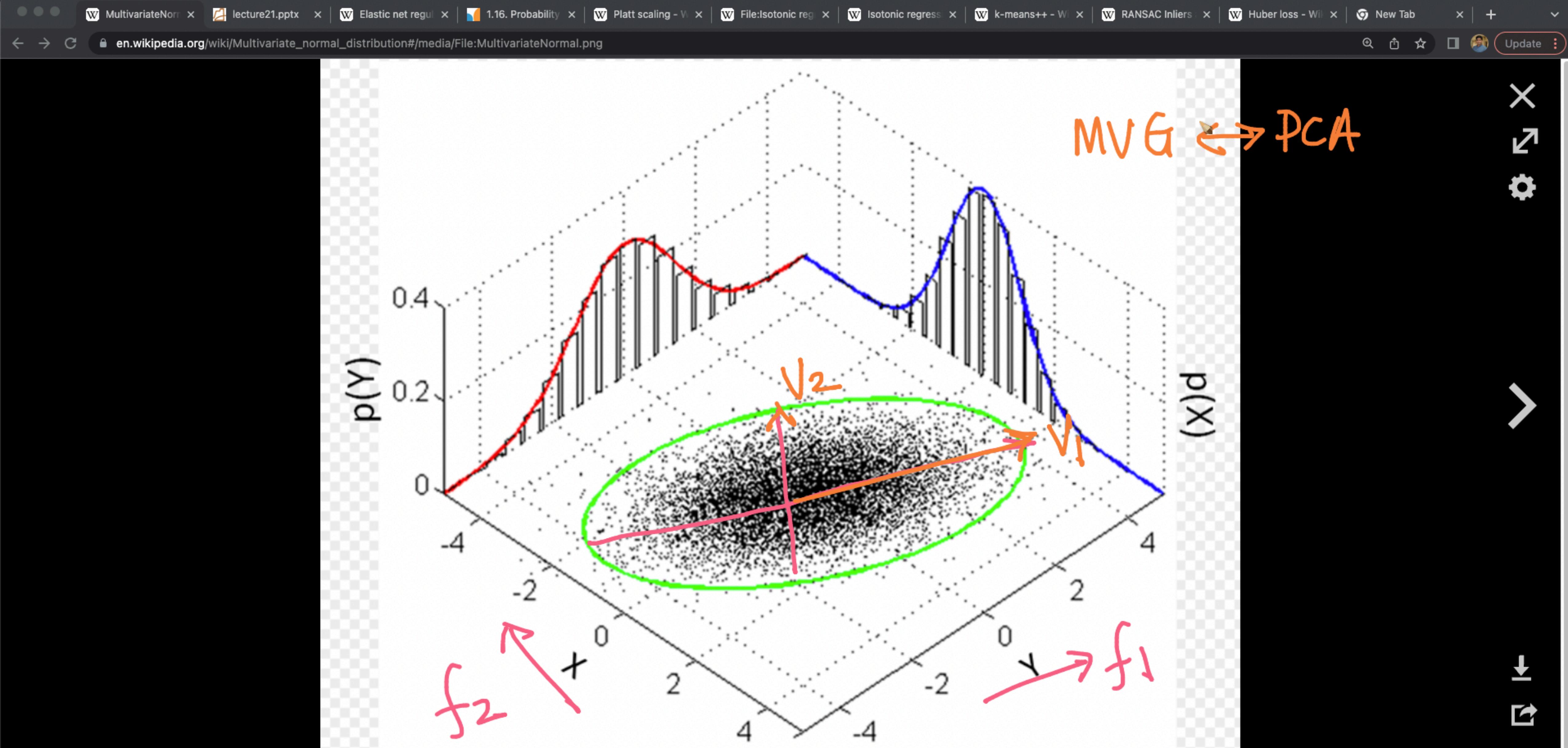
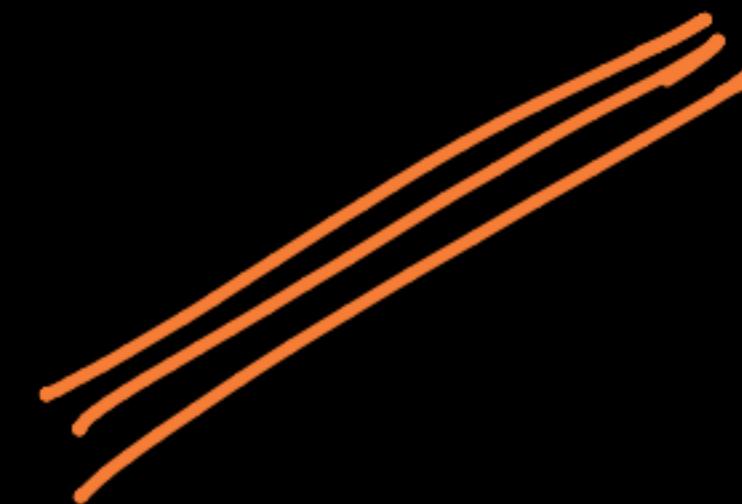


Illustration of a multivariate gaussian distribution and its marginals. Matlab code provided below.

More details



x_1, x_2, \dots, x_n

$$x \sim N(\mu, \sigma)$$

1D-Gaussian:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$



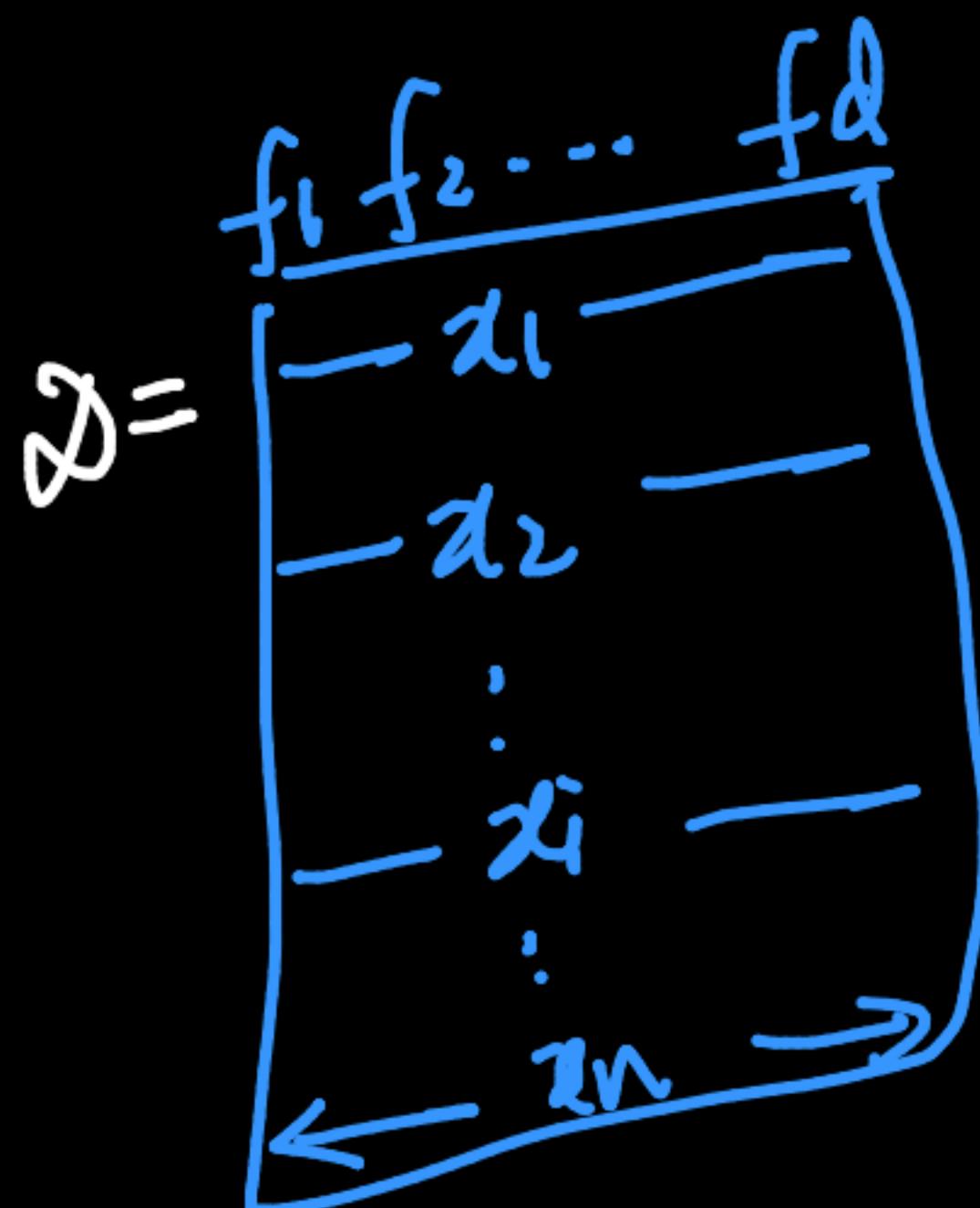
$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

d-dim Gaussian:

dxd symm - Matrix

params

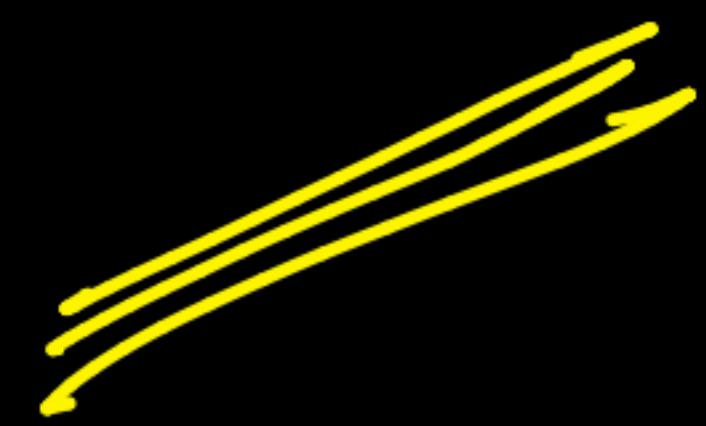
$$\mathcal{D} : \underbrace{x_1, x_2, \dots, x_i, \dots, x_n}_{\text{d-dim vec}} : X \sim N_d \left(\underline{\mu}_d, \underline{\Sigma}_{d \times d} \right)$$



$$+ x_i \in \mathbb{R}^d$$

$$\left\{ \begin{array}{l} \hat{\mu}_d = \frac{1}{n} \sum_{i=1}^n x_i \\ \text{d-dim vec} \end{array} \right.$$

$$\left\{ \begin{array}{l} \hat{\Sigma}_{ij} = \text{cov}(f_i, f_j) \\ \text{PCA (earlier)} \end{array} \right.$$



$$\mu_d^j, \Sigma_{d \times d}^j \quad \sum_j = k$$

Find $\underset{=}{k}$ -Gaussians

Given: $\mathcal{D} = \{\underline{x}_i\}_{i=1}^n \quad \underline{x}_i \in \mathbb{R}^d$

$\leftarrow \{ \begin{matrix} j: \text{index of the Gaussian} \\ j; l \rightarrow k \end{matrix} \right.$

Generative Methods

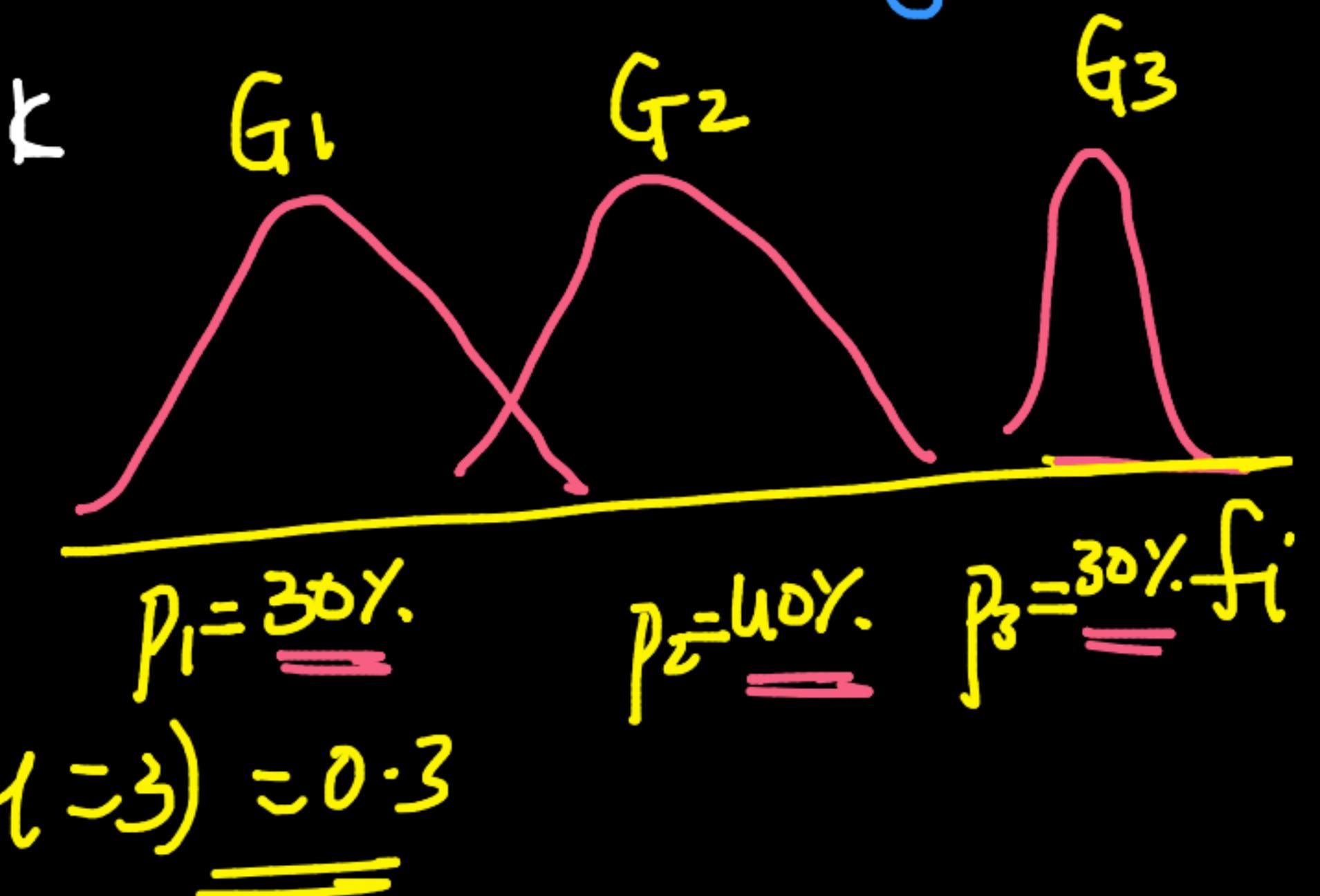
GMM

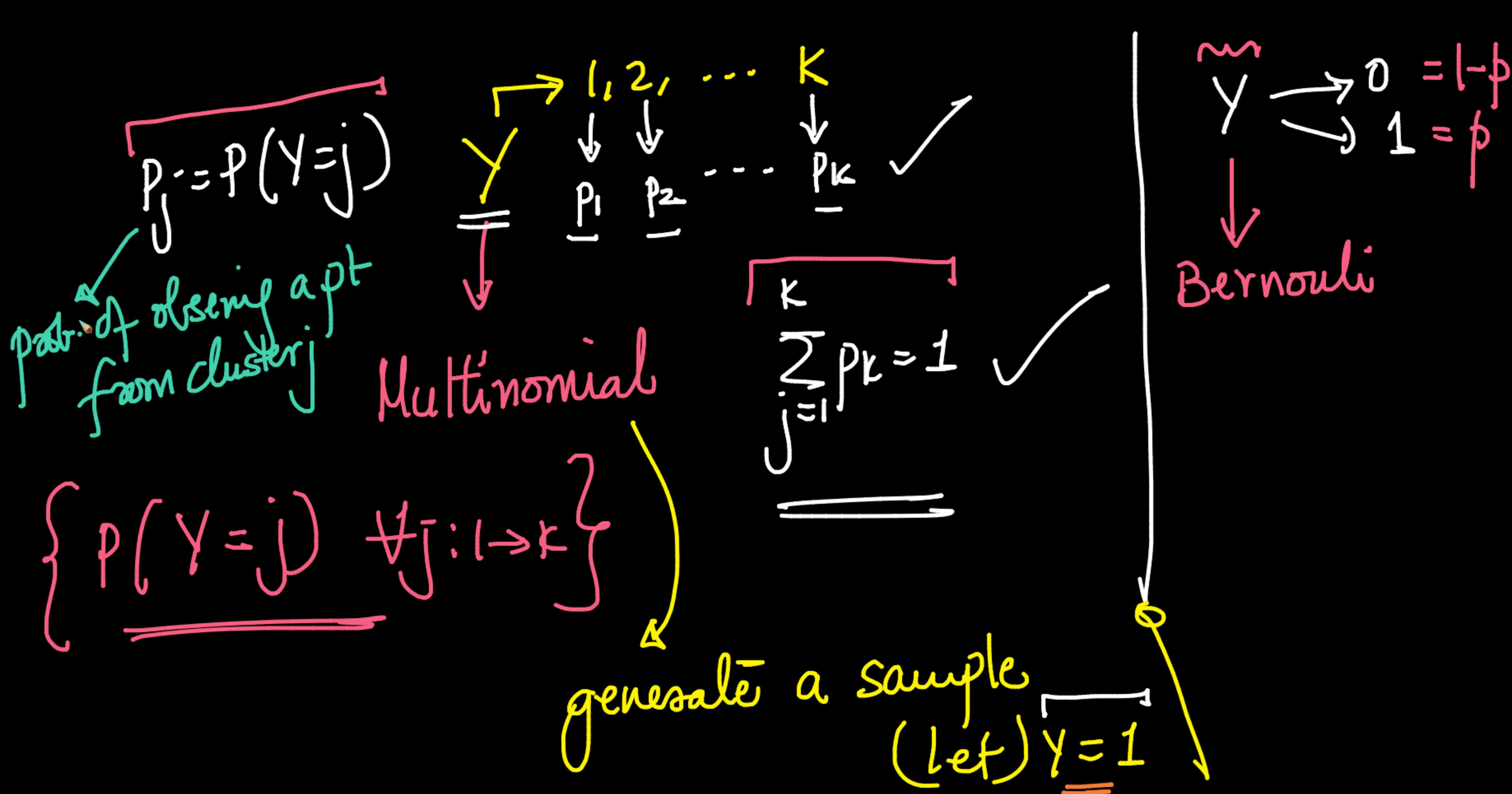
Assumptions on how the data was generated

Generative process
underlying GMM

* a $\rightarrow P(Y=j) \quad j: 1 \rightarrow k$

$$\left\{ \begin{array}{l} P(Y=1) = 0.3 \\ P(Y=2) = 0.4 \\ P(Y=3) = 0.3 \end{array} \right.$$





6

⊕ $\mu_d^j, \Sigma_{d \times d}^j \quad \forall j = 1 \rightarrow K$

↳ params of each underlying Gaussian

$\therefore Y = 1$

Sample from $N_d(\mu_d^j, \Sigma_{d \times d}^j)$

↳ one data point

Repeat steps a & b many times

↓

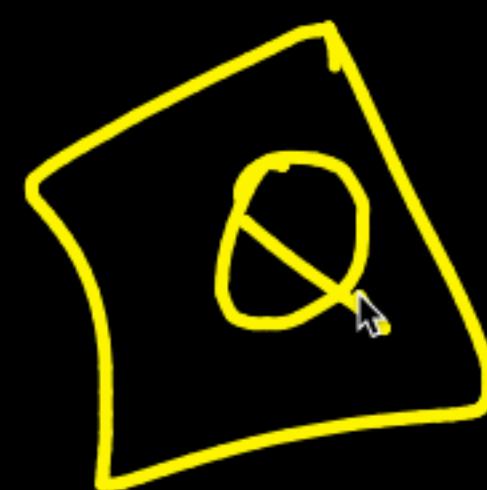
you can obtain data very "similar" to A

Discriminative - methods

all
techniques
fill now

$x_i \rightarrow \text{cl}_0 \text{ or } \text{cl}_1$

$x_i \rightarrow \text{cluster}_0 \text{ or } \text{cluster}_1$



$$\left\{ \underbrace{P(Y=j)}_{j=1} = p_j \right.$$

$$\sum_{j=1}^k p_j = 1$$

$\underbrace{1, 2, \dots, j, \dots, k}$

① $P(Y_i=j | \underline{x_i})$

②

$$\sum_{j=1}^k P(Y_i=j | \underline{x_i}) = 1$$

Q

GMMs

Strong assumption ①

Seldom used

- Generative model

- Optimization (MLE: later)

- EM-algo (alt. to GD)

K - Underlying Gaussians

hard
GMM \approx K-means

②

Recap:

- Compute $\mu_d, \Sigma_{d \times d}$ given $\mathcal{D}: x_i \in \mathbb{R}^d$
- Generative process underlying GMM
(assumption)

→

find

$$\underbrace{P(Y=j)}_{\text{all params: } \Theta} \forall j = 1 \rightarrow K$$

and $\mu_d, \Sigma_{d \times d} \forall j: 1 \rightarrow K$

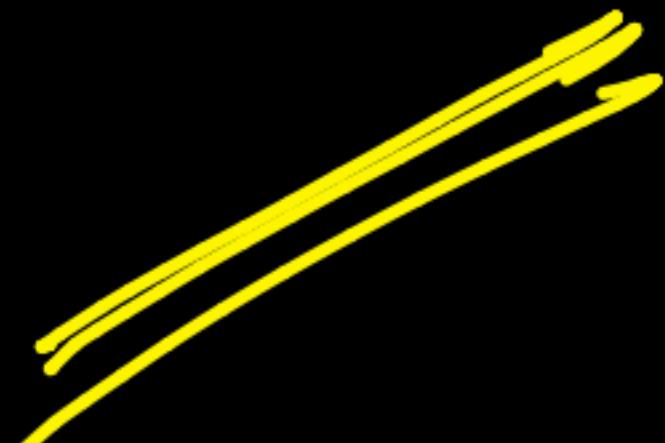
$$\sum_{j=1}^k P(Y=j) = 1$$

using

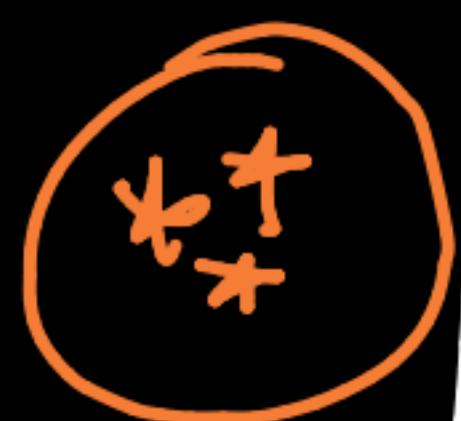
$$\mathcal{D}: \{x_i\}_{i=1}^n \quad x_i \in \mathbb{R}^d$$

$$\rightarrow \Theta = \left[\underbrace{P(Y=j) \forall j}_{\text{all params: } \Theta}, \underbrace{\mu_d, \Sigma_{d \times d}}_{\text{all params: } \Theta} \forall j \right]$$

all params: Θ
(let)



Given: $\text{observed } \mathcal{D} := \{x_i\}_{i=1}^n$



Find θ s.t. the probability of generating \mathcal{D} is maximal

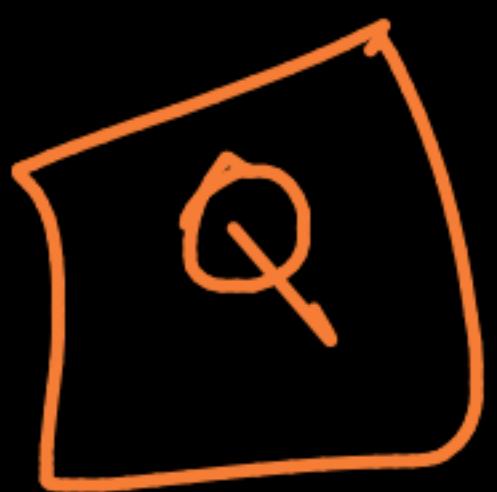


→ using the 2-step generative process

Max-Likelihood
Estimation

↓
probability meets optimization

extensively later



ha
GMM
generative
method

\approx

K-means
(discriminative)

\uparrow

GMM

\approx

soft K-means



$$\underset{\theta}{\operatorname{Max}} \quad P(D)$$

{
 [
 gen salive = process
 underly ing GMM

$$\Rightarrow \underset{\theta}{\operatorname{Max}} \quad P(x_1, x_2, \dots, x_n)$$

$$D = \{x_i\}_{i=1}^n$$

$$\Rightarrow \underset{\theta}{\operatorname{Max}} \quad \prod_{i=1}^n P(x_i)$$

↑
post. of generativ x_i

each x_i is
indep of one
another

Max $\prod_{i=1}^n P(x_i)$ → Post of observing x_i Generative process:

$= \max_{\theta} \prod_{i=1}^n \sum_{j=1}^K P(y_j, x_i)$

a) one j using P_j 's

b) generative data from $N(\mu_d, \Sigma_d)$

Post of observing x_i being generated by j th Gaussian

$$= \max_{\theta} \prod_{i=1}^n \sum_{j=1}^K P(y_j, x_i) \xrightarrow{P(y_j \wedge x_i)}$$

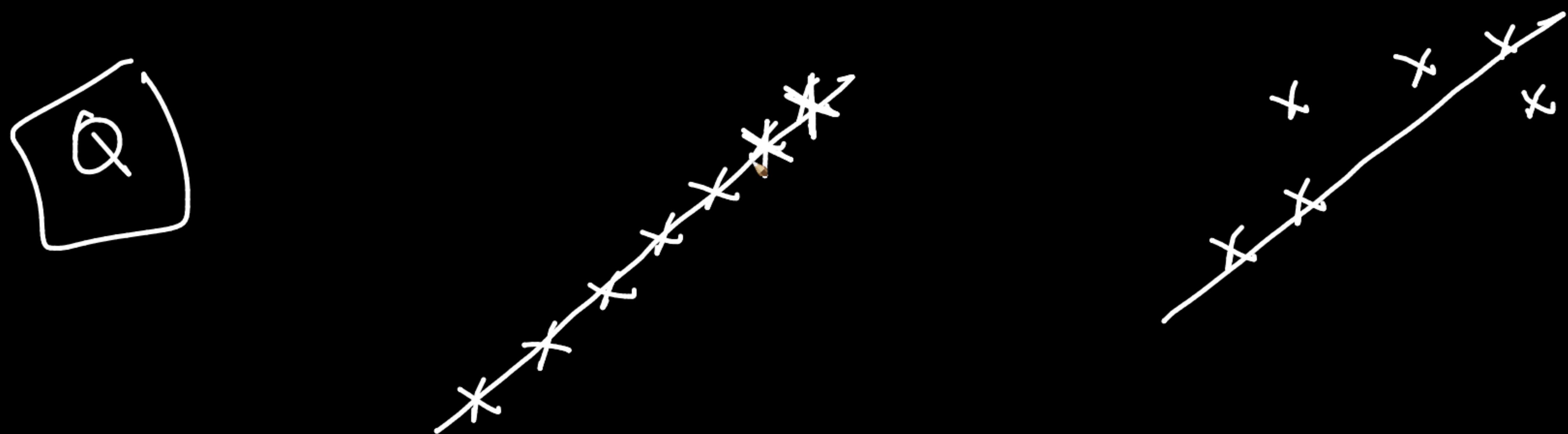
conditional-path

$$= \max_{\theta} \prod_{i=1}^n \sum_{j=1}^K \left(P(x_i | y_j) P(y_j) \right)$$

↳ class priors

Params

$\theta = [P(Y=j); \mu_d, \Sigma_{dxd}; j:1 \rightarrow K]$ likelihood of generating x_i using j^{th} Gaussian



$$\underset{\theta}{\operatorname{Max}} \prod_{i=1}^n \sum_{j=1}^K P(x_i | \tilde{y}_j) = P(y_j)$$

$$\theta = \underline{P(y=j)} ; \underline{\mu^j, \Sigma^j}$$

$$p_j = P(y_j)$$

plug in the PDF of
 $N(\mu_d, \Sigma_{dxd})$

PMF of α
 multinomial r.v

The spherical normal distribution can be characterised as the unique distribution where components are independent in any orthogonal coordinate system.^{[3][4]}

Density function [edit]

Non-degenerate case [edit]

The multivariate normal distribution is said to be "non-degenerate" when the symmetric covariance matrix Σ is positive definite. In this case the distribution has density^[5]

$$f_{\mathbf{x}}(x_1, \dots, x_k) = \frac{\exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}))}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}}$$

where \mathbf{x} is a real k -dimensional column vector and $|\boldsymbol{\Sigma}| \equiv \det \boldsymbol{\Sigma}$ is the determinant of $\boldsymbol{\Sigma}$, also known as the generalized variance. The equation above reduces to that of the univariate normal distribution if $\boldsymbol{\Sigma}$ is a 1×1 matrix (i.e. a single real number).

The circularly symmetric version of the complex normal distribution has a slightly different form.

Each iso-density locus — the locus of points in k -dimensional space each of which gives the same particular value of the density — is an ellipse or its higher-dimensional generalization; hence the multivariate normal is a special case of the elliptical distributions.

The quantity $\sqrt{(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})}$ is known as the Mahalanobis distance, which represents the distance of the test point \mathbf{x} from the mean $\boldsymbol{\mu}$. Note that in the case when $k = 1$, the distribution reduces to a univariate normal distribution and the Mahalanobis distance reduces to the absolute value of the standard score. See also Interval below.

Bivariate case [edit]

In the 2-dimensional nonsingular case ($k = \text{rank}(\boldsymbol{\Sigma}) = 2$), the probability density function of a vector $[XY]'$ is:

$$f(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y \sqrt{1 - \left[\frac{1}{\sigma_X^2} \left(\frac{(x - \mu_X)^2}{\sigma_X^2} + \frac{(y - \mu_Y)^2}{\sigma_Y^2} \right) \right]}}$$

Multivariate Normal Distribution
Bivariate normal joint density

$$\underset{\theta}{\operatorname{Max}} \prod_{i=1}^n \sum_{j=1}^K P(x_i | y_j) \cdot \underbrace{P(y_j)}_{\substack{\text{One } A \text{ like} \\ \text{params}}}$$

~~Complex~~

- complex to solve ✓
- GD - complex
- multiple local minima

PDF · A MDG
 μ_d^j, Σ_d^j
 $\Sigma_w, [\Sigma]$...

'Hack' to solve the optimzn w/o using GD ✓

↳ Expectation - Maximization :

$$\theta = [P_i^{\text{K}}, \mu_d^{\text{KxK}}, \Sigma_{d \times d}^{\text{KxK}}, H_j^{\text{J=1 to K}}]$$

Code-idea!

→ { E: update; fix ; M: fix ; update }

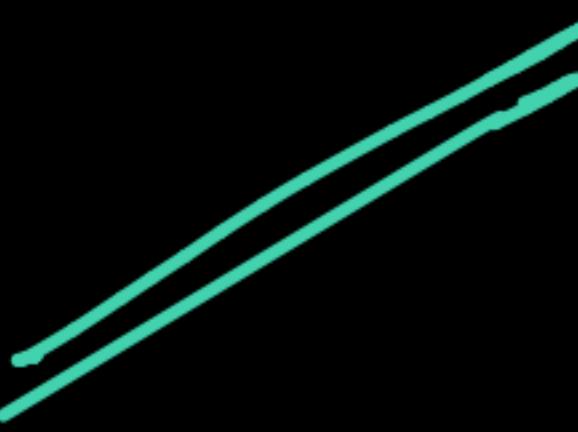
Computationally
much cheaper
than GD

EM: coordinate ascent $\underset{\theta_1, \theta_2}{\text{Max}} f(\theta_1, \theta_2)$

$\theta_2 = \text{param}_2$

$\theta_1 = \text{param}_1 = \theta_1$

The diagram illustrates the iterative nature of the EM algorithm. It starts at a point on the parameter space (indicated by a small red dot), moves vertically upwards (along the θ_2 axis) to a local maximum (the 'optima' point marked with a blue circle), and then moves horizontally to the left (along the θ_1 axis) to the next iteration point.



Recap:

n-dim

→ Gaussian Mixture-Model →

θ ; p_j ; μ_d ; $\Sigma_{d \times d}$
 $j; l \rightarrow k$

→ Generative process $\rightarrow p_j, \cdot$
 $\mu_d, \Sigma_{d \times d}$

→ Max $P(\theta)$

break: 10:35

GMM

Given: \underline{x}_i 's $i=1 \rightarrow n$; $\underline{\mu}, \underline{\Sigma}$ (very similar to k-Means)

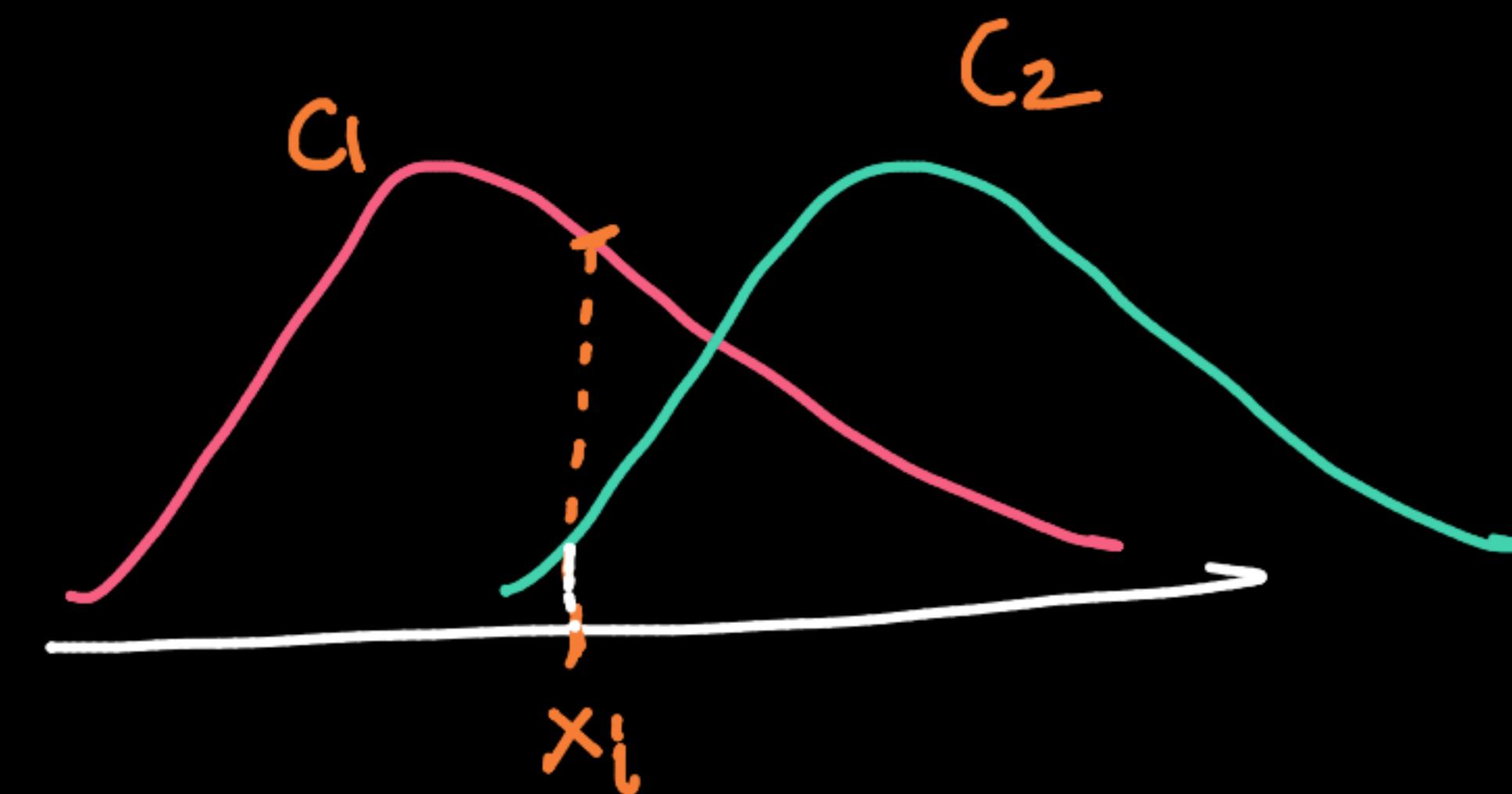
① Exp: soft assignment

for each \underline{x}_i , we compute the prob it belongs to jth cluster (PDF) $P(\underline{x}_i \in j^{\text{th}})$ using PDF + normalize

② Max:

(update)

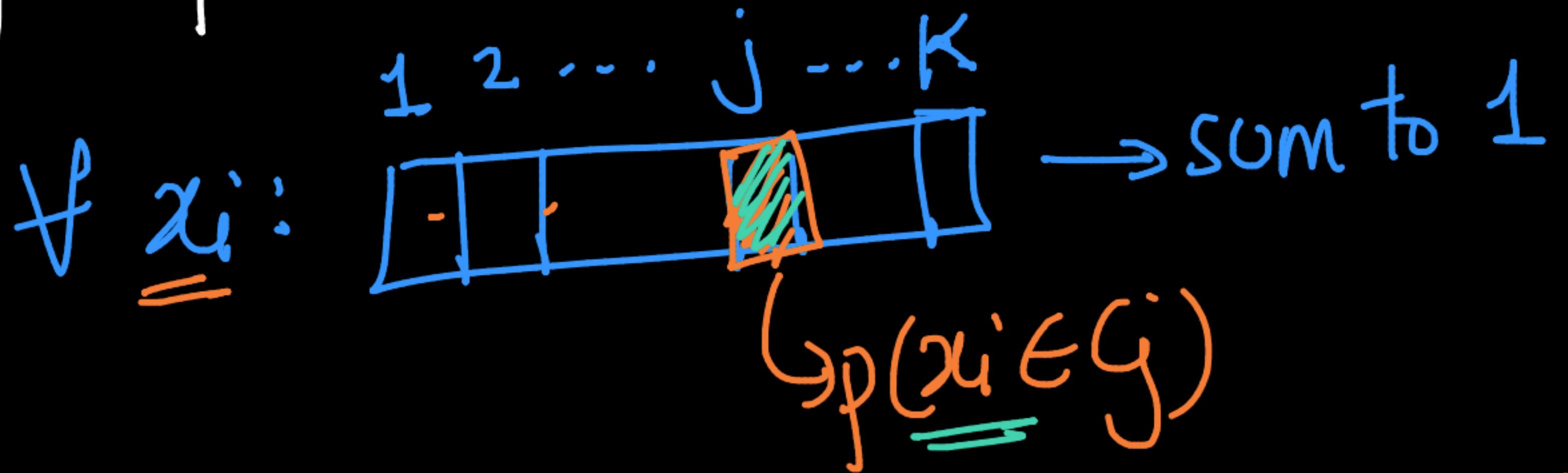
Re-estimate Gaussian - params for all Gaussians
 $\underline{\mu}$ & $\underline{\Sigma}$ using the weighted scheme
 update $\underline{\mu}$ & $\underline{\Sigma}$



$\frac{p_1}{p_1 + p_2}$ $f(x_i \in C_1)$: PDF = $\underline{p_1}$

$\frac{p_2}{p_1 + p_2}$ $f(x_i \in C_2)$: PDF = $\underline{p_2}$

@ end of exp-step:-

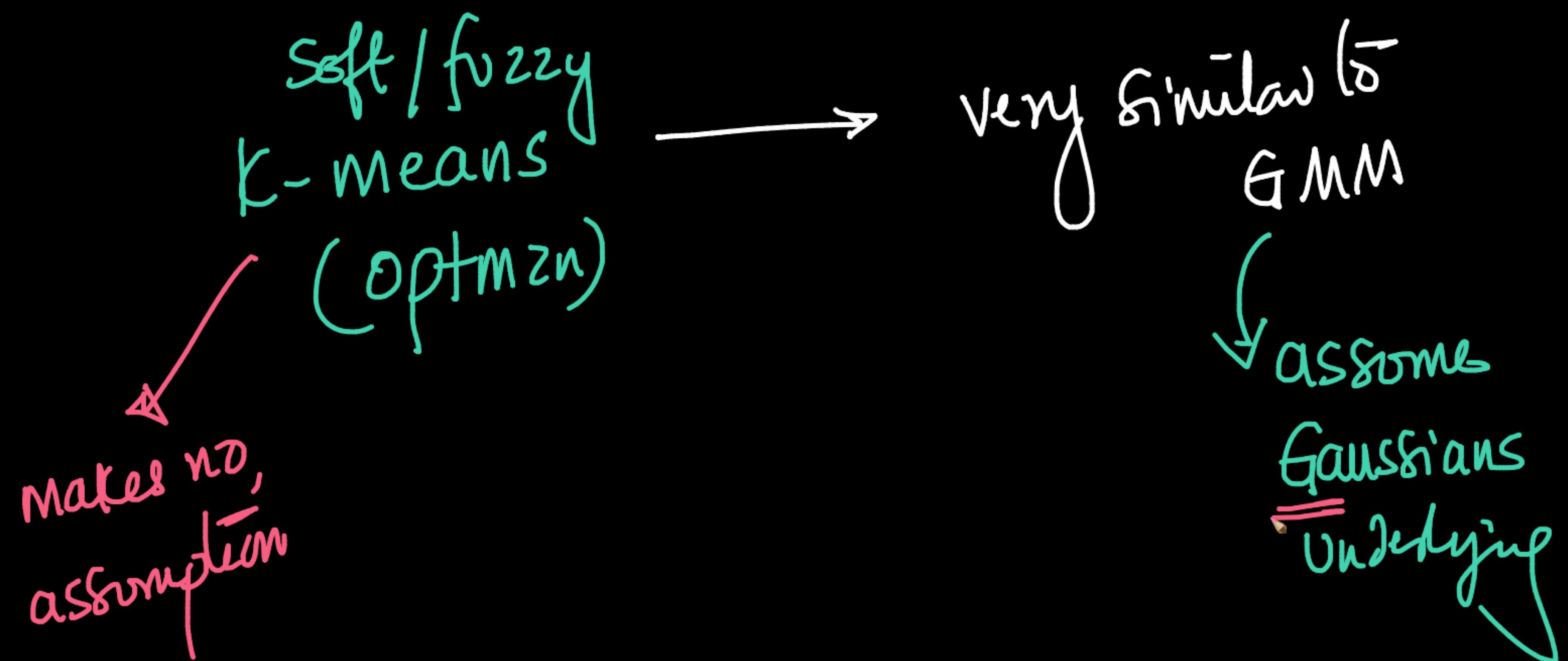


✓ In Max-step:

Same as
Centroid

μ^j = weighted average of each x_i
 $\hookrightarrow p(x_i | G_j)$

Σ^j : weighted Co-var =



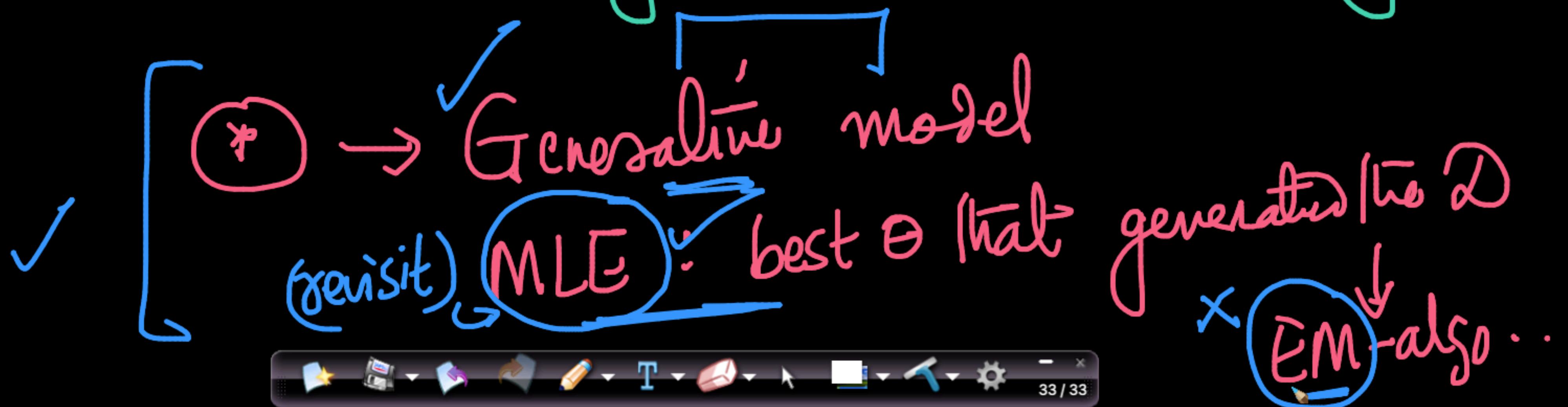
~~GMM~~: \rightarrow soft-clustering . . .

\rightarrow less used in practice

\therefore hard GMM = k-means (popular)



\rightarrow Strong assumption about Gaussian mixture of underlying the data

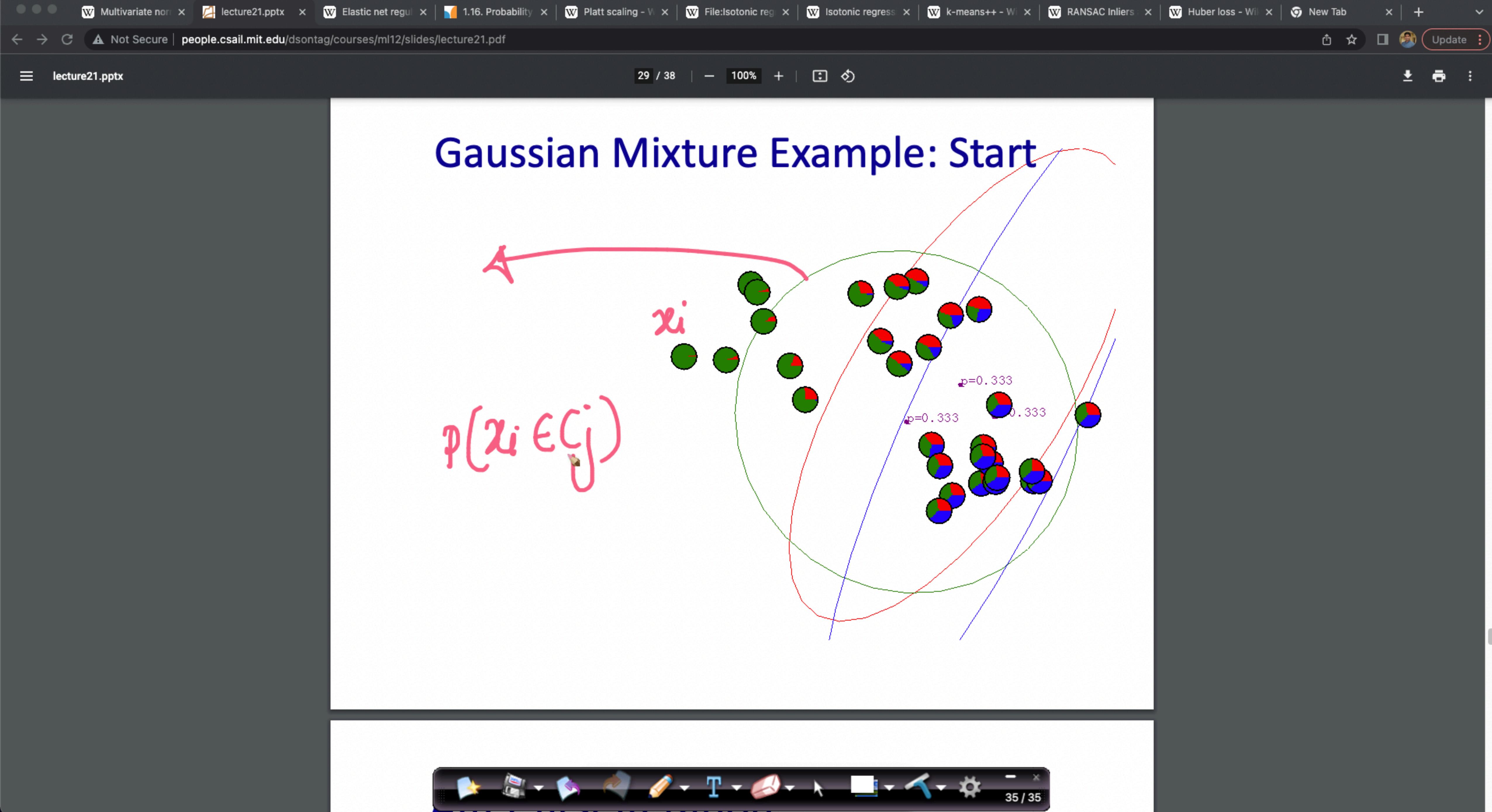


MLE; MAP; Bayesian-opt



linear-regression

Post-sead
References
logistic-reg



Multivariate norm x lecture21.pptx x Elastic net regul x 1.16. Probability x Platt scaling - W x File:Isotonic reg x Isotonic regress x k-means++ - Wi x RANSAC Inliers x Huber loss - Wi x New Tab x +

Not Secure | people.csail.mit.edu/dsontag/courses/ml12/slides/lecture21.pdf

lecture21.pptx 29 / 38 - 100% + ☰

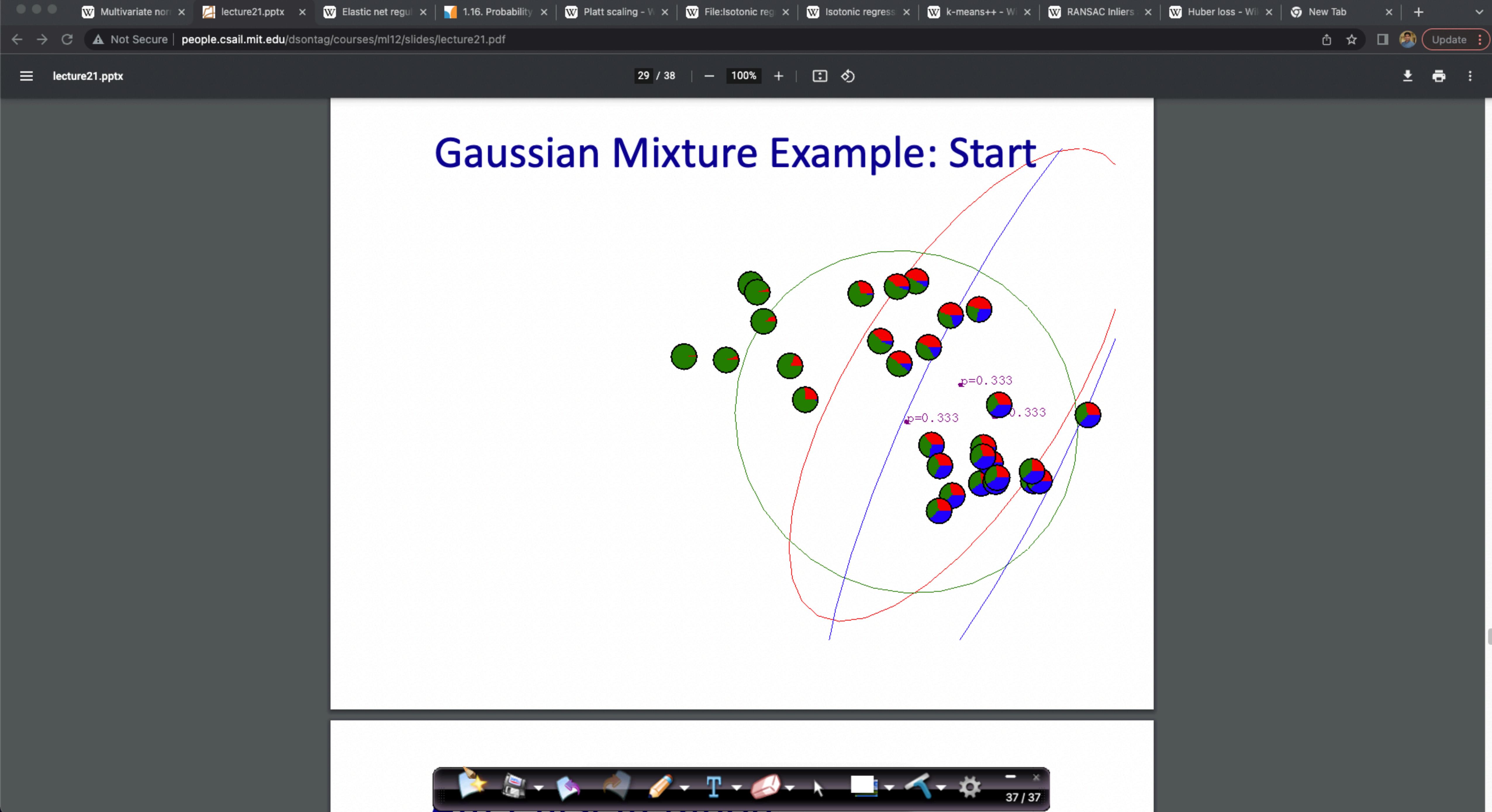
Gaussian Mixture Example: Start

random ✓ $\mu_d^j, \Sigma_{d \times d}^j = j: 1 \rightarrow 3$

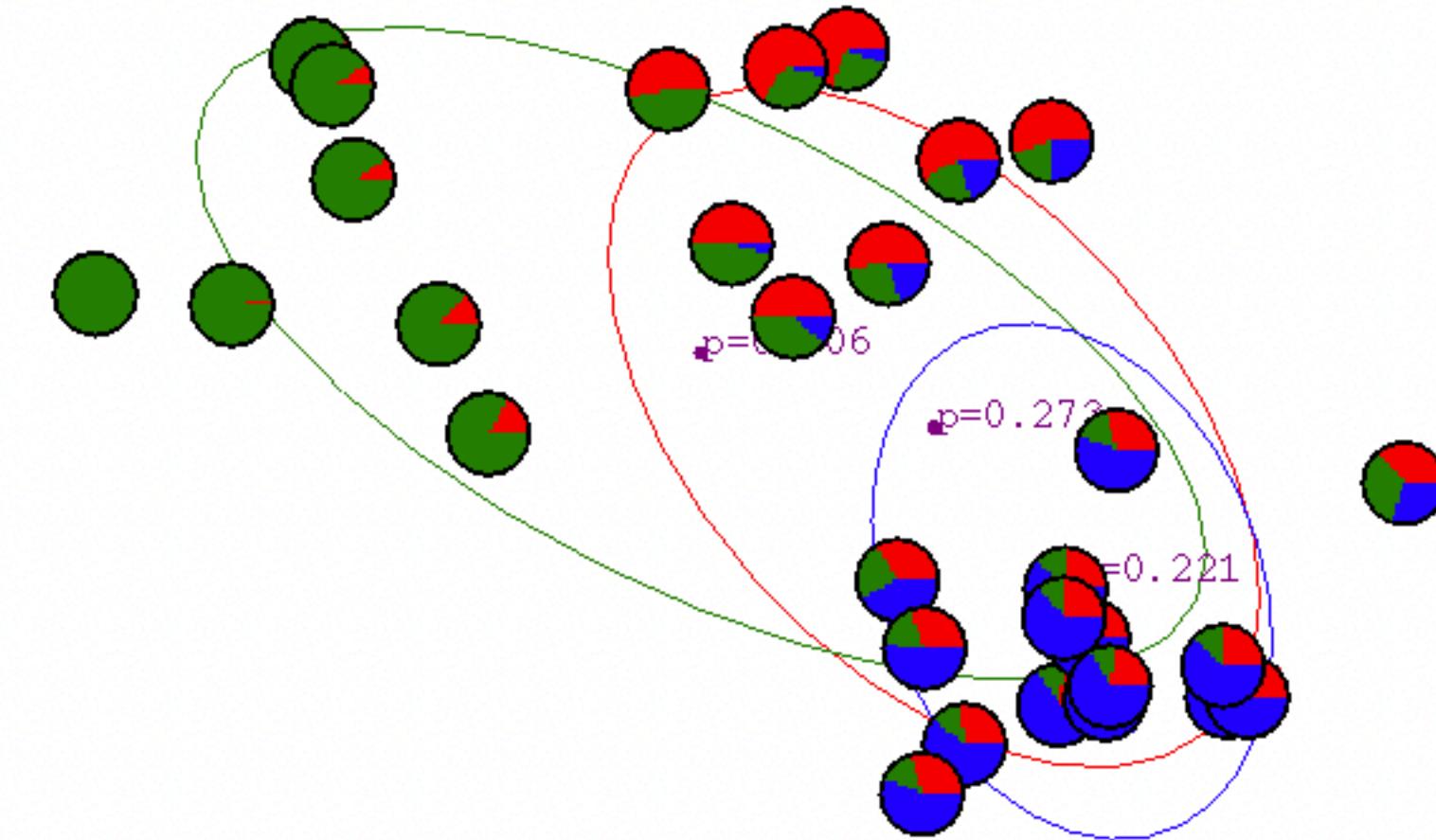
Exp: soft assign

Max: recompute μ_j, Σ_j

36 / 36



After first iteration



Multivariate norm x lecture21.pptx x Elastic net regul x 1.16. Probability x Platt scaling - W x File:Isotonic reg x Isotonic regress x k-means++ - Wi x RANSAC Inliers x Huber loss - Wi x New Tab x +

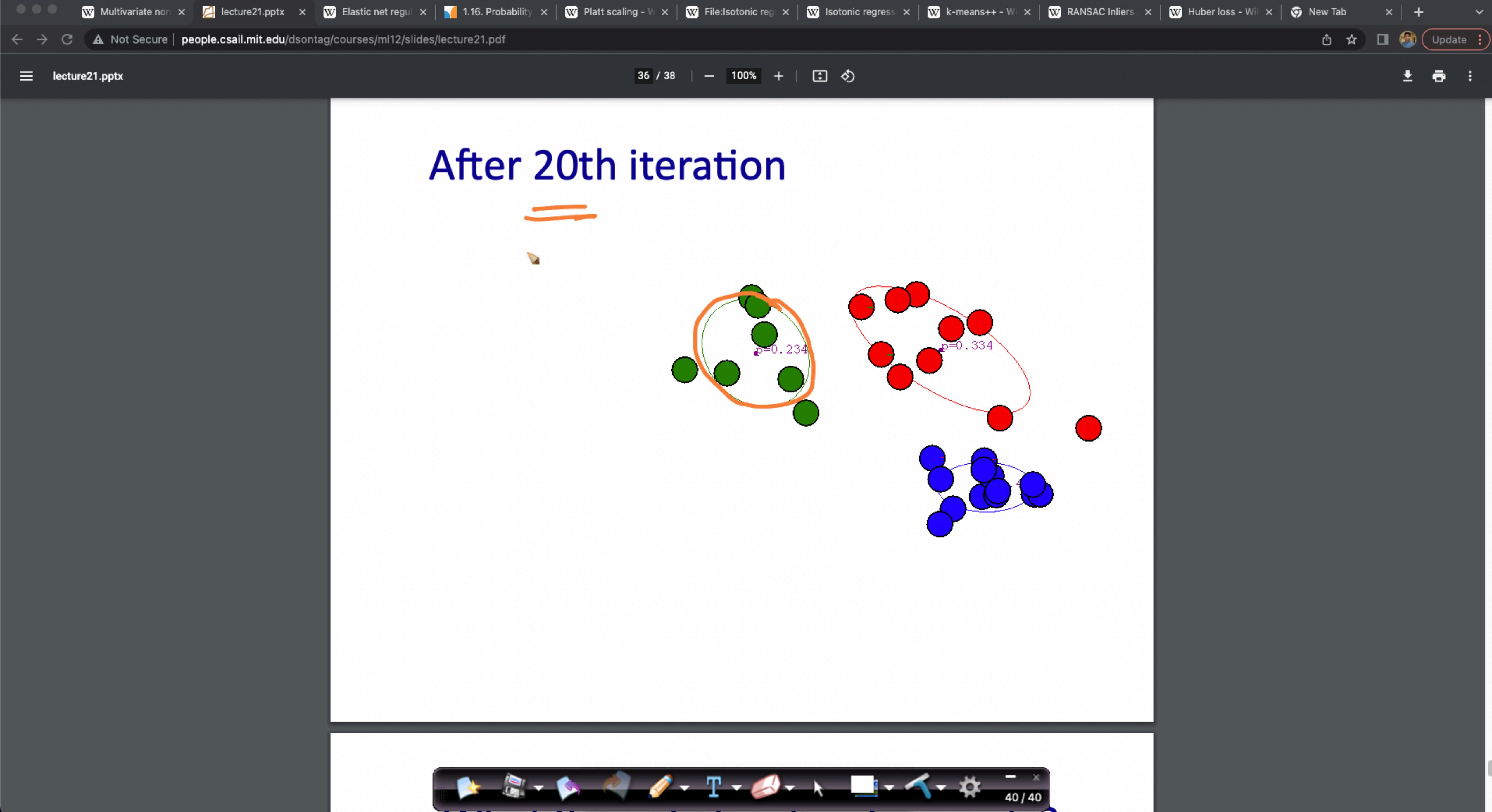
Not Secure | people.csail.mit.edu/dsontag/courses/ml12/slides/lecture21.pdf

34 / 38 | - 100% + | ☰ ⚡

lecture21.pptx

After 5th iteration

39 / 39



Suggestion

— MLE / optim2n } revisit

— EM } revise

Misc. Topics

L₂-reg:

Ridge-regression

= L₂-reg + linear reg

Tikhonov reg ...

log-loss
sq. loss
hinge loss
huber loss

$$\min_{w_j} \text{loss} + \lambda \sum_{j=1}^d w_j^2$$

$$\sum_{j=1}^d w_j^2 \rightarrow L_2\text{-norm}$$

weights of useless features to become small

Why $L_2\text{-reg}$ in Linear-regression ?
↓
reduce overfitting

L₁-reg or Lasso

$$\min_{w_j} \text{loss} + \lambda \sum_{j=1}^d |w_j|$$

↑
sq. loss

L₁-norm of w

Is there a
problem here?

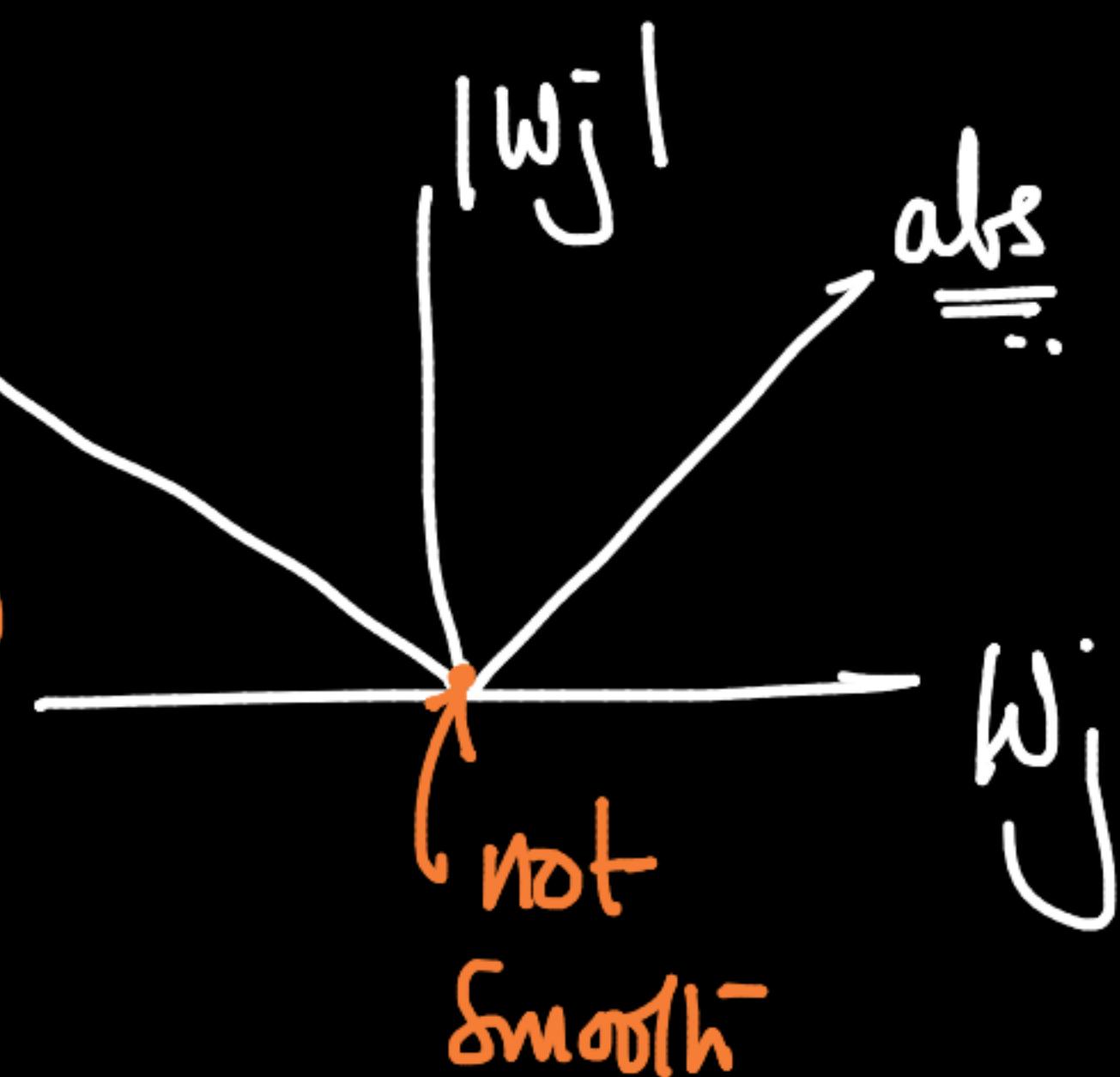


$$L_1\text{-norm}_{(w)} = \sum_{j=1}^d |w_j|$$

$$\min_{\omega_j} \mathcal{L} = \sum_{i=1}^n (y_i - (\omega^T x_i + \omega_0))^2 + \lambda \sum_{j=1}^d |\omega_j|$$

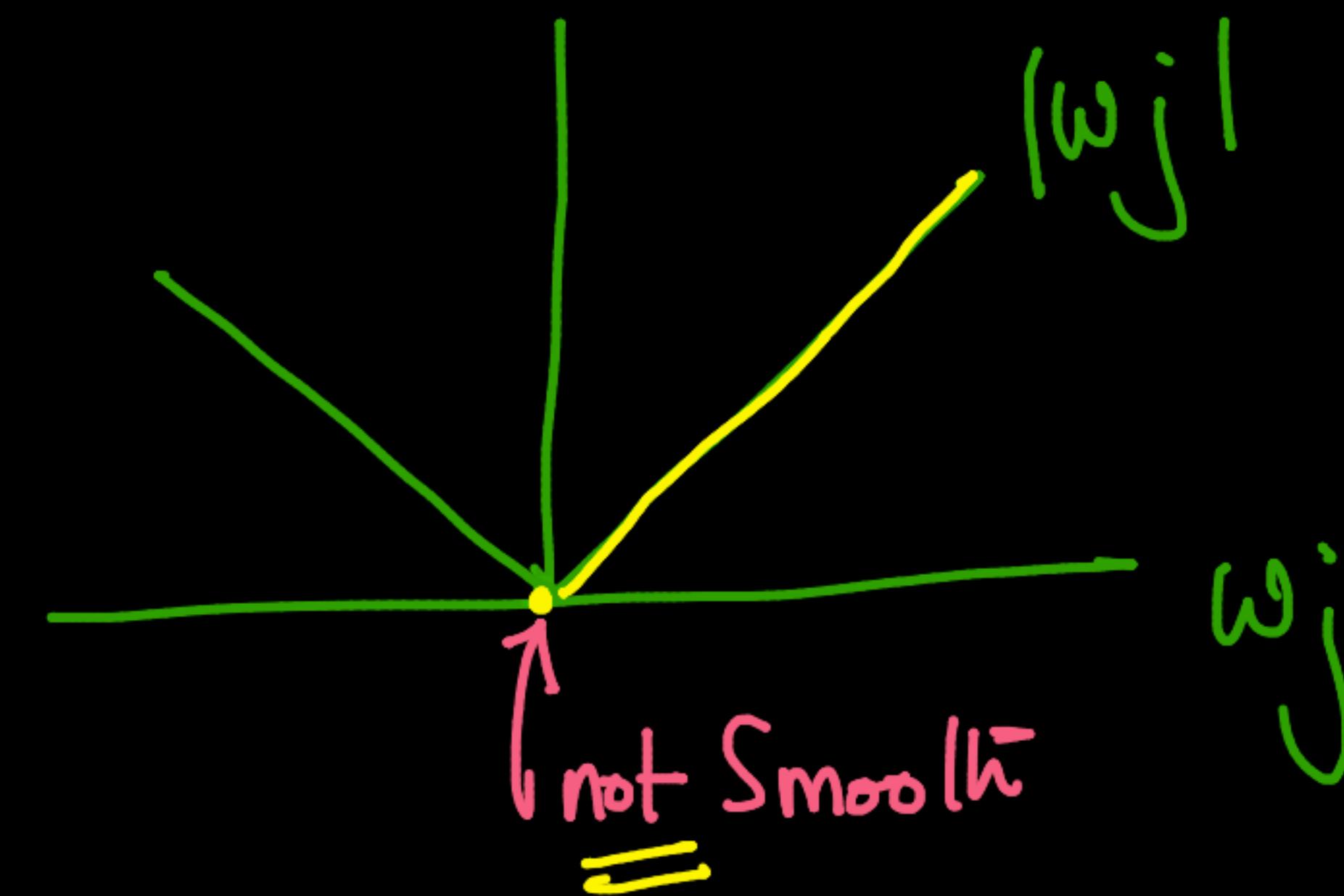
HINT: GD

$$\left\{ \begin{array}{l} \frac{\partial \mathcal{L}}{\partial \omega_j} \\ \text{is not defined} \end{array} \right|_{\omega_j=0}$$



fix \bar{w} :

$\text{der-abs}(\bar{w}_j)$



$$\frac{d|w_j|}{dw_j} = \begin{cases} +1 & \text{if } w_j > 0 \\ -1 & \text{if } w_j < 0 \\ 0 & \text{if } w_j = 0 \end{cases}$$

hack

$$\frac{d [w_j]}{dw_j} @ w_j = 0 \\ \downarrow = 0 \text{ (why is it } \underline{\underline{\text{ok}}}\text{?})$$

↳ ∵ $w_j = 0$ means that
 f_j is useless
⇒ we can keep
 $w_j @ 0$

GD:

→ init all w_j 's randomly + non-zero

→ update w_j 's

if some $w_j = 0 \Rightarrow f_j$ is
useless

$w_j^{\text{new}} = w_j^0 - \eta \frac{\partial L}{\partial w_j}$

$w_j^0 \leftarrow \frac{\partial L}{\partial w_j}$ is not defined
(OK)

~~L₁-reg:~~

→ results in a Sparse soln \Rightarrow w is sparse

Why?
~~if~~

[all less useful features
weights become $= 0$]

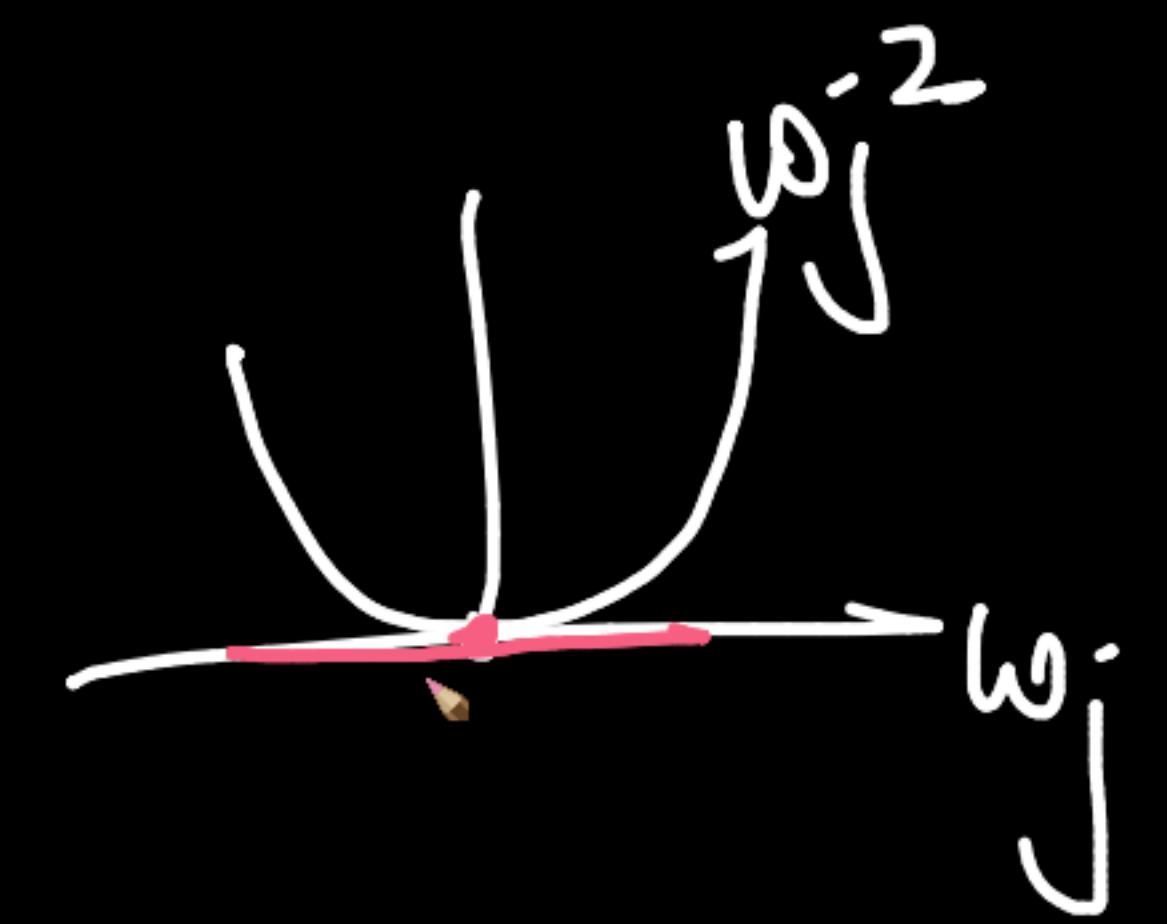
L₂-reg: \rightarrow all less useful features weights
become very small & close to 0

but often not exactly
equal to 0

$$\frac{d|\omega_j|}{d\omega_j} \Big|_{\omega_j=0} = 0$$

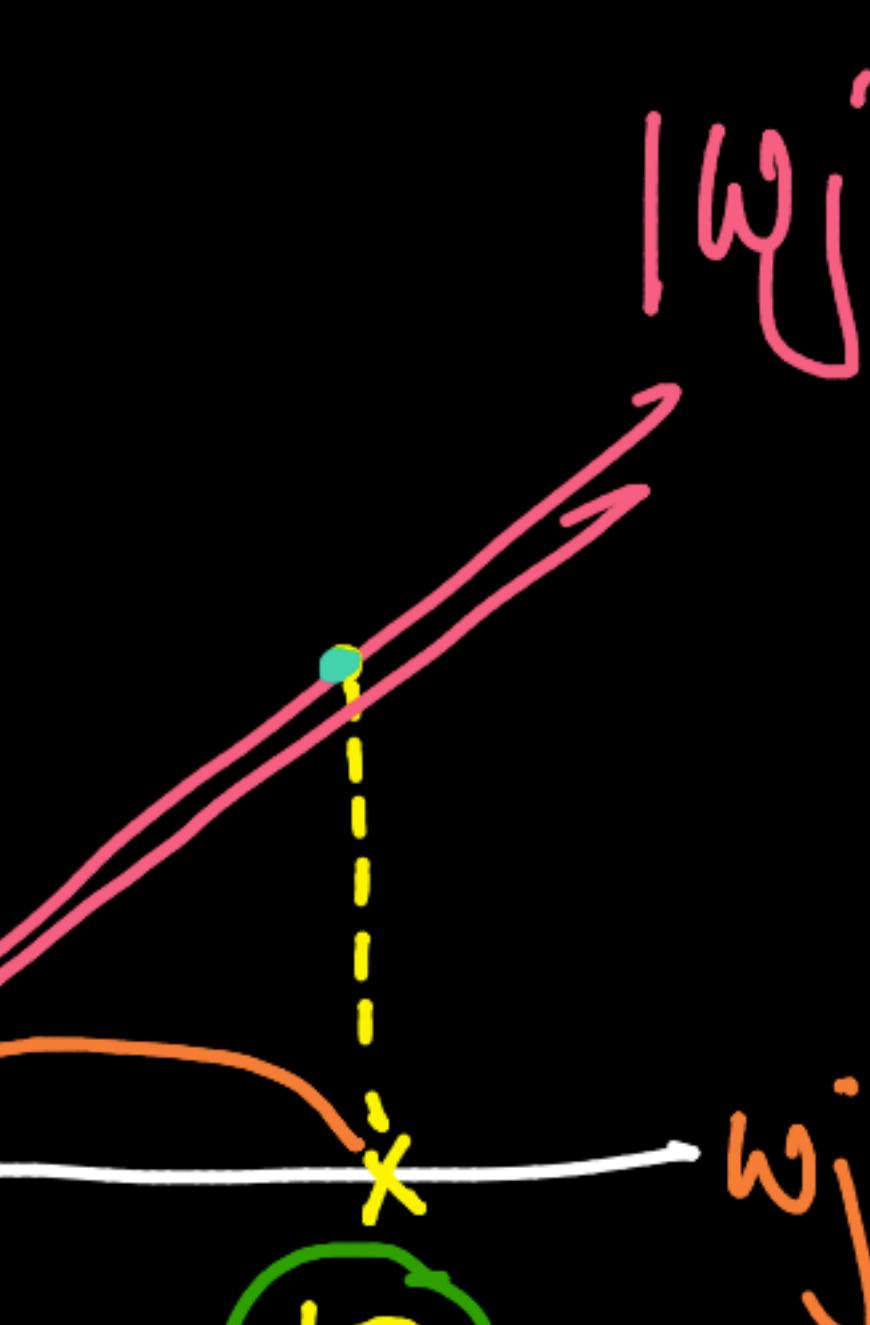


$$\frac{d\omega_j^2}{d\omega_j} \Big|_{\omega_j=0} = 0$$



$$w_j = 0$$

$L_1\text{-reg}$



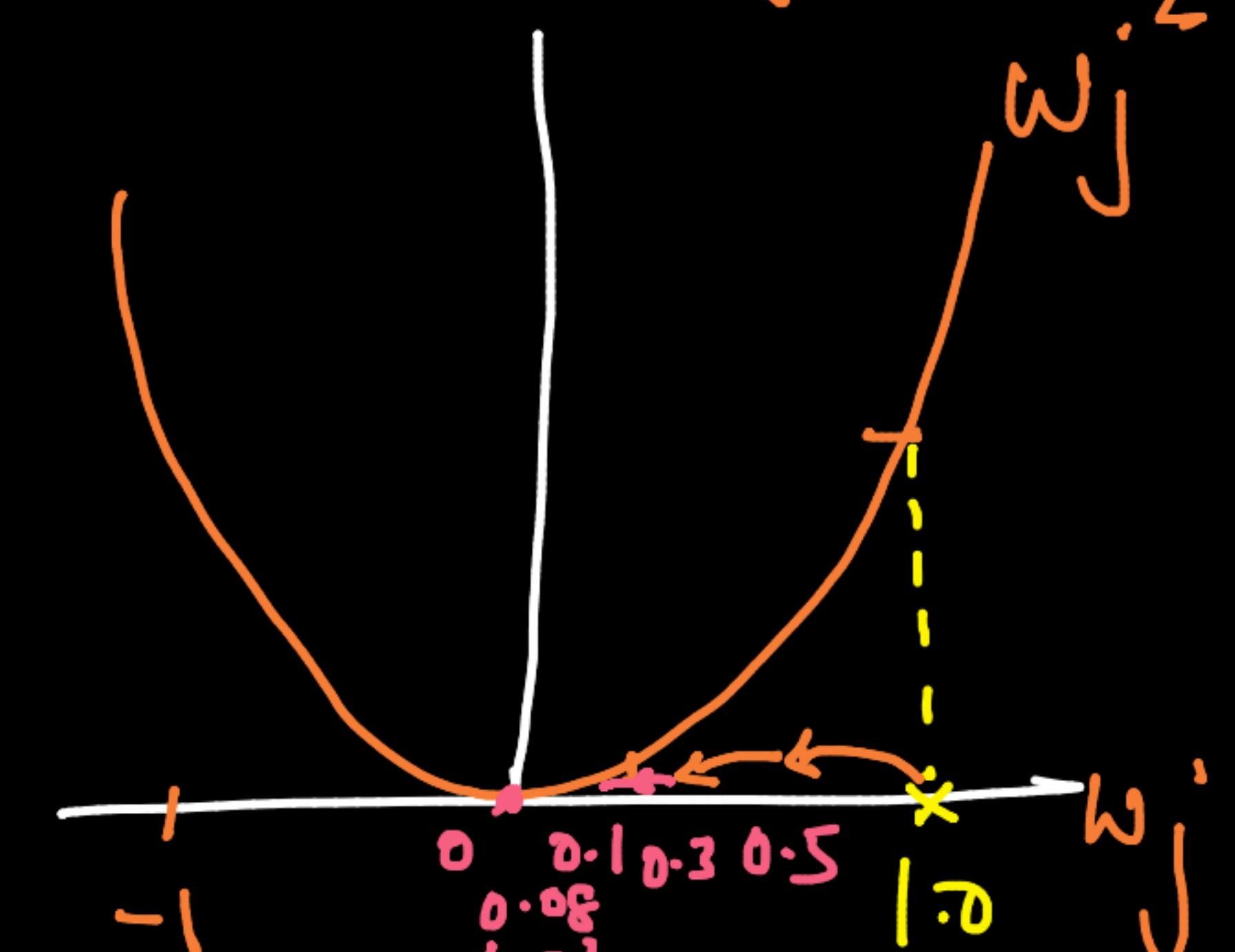
$$w_j^{\text{new}} = w_j^{\text{old}} - \eta \frac{\partial L}{\partial w_j} + 1$$

0.5 (let)

lets assume f_i is a useless feature

geometric

$L_2\text{-reg}$



$$\frac{\partial L}{\partial w_j} = 2w_j$$

In L₂-reg as we approach '0'
The gradient also reduces

$$w_j^{\text{new}} = w_j^{\text{old}} - \eta \frac{\partial \tilde{L}}{\partial w_j}$$

$L_1\text{-reg} \rightarrow$ better Sparsify than $L_2\text{-reg}$

disadv:

i

dim #data points
 $d > n$ (genomics)

$\hookrightarrow L_1\text{-reg}$ will have utmost n non-zero features

even if $\geq n$ features are useful

② highly correlated features

↳ one feat non-zero
↳ all other feat: zero ✓

↓

using weights fw
F-I
(mess-up)

W Multivariate norm x | A lecture21.pptx x | W Elastic net regul x | 1.16. Probability x | W Platt scaling - W x | W File:Isotonic reg x | W Isotonic regress x | W k-means++ - Wi x | W RANSAC Inliers x | W Huber loss - Wi x | New Tab x | +

Not Secure | people.csail.mit.edu/dsontag/courses/ml12/slides/lecture21.pdf

34 / 38 | - 100% + | ☰ ⚡

lecture21.pptx

After 5th iteration

55 / 55