

# Classification & Regression

## (classical ML)

# List

KNN

Linear Regr

Logistic Regr

Naive Bayes

Bagging : RF

Adaboost

Boosting : GBDT; CatBoost

LightGBM

Stacking, Cascading

SVM

SVR

L<sub>1</sub>

kernel

SMOTE

extra Trees

MLE, MAPE

Regularization : L<sub>1</sub>, L<sub>2</sub>, ElasticNet

Calibration : Platt-scaling, Isotonic reg

RANSAC

bias-variance

Loss fns

Misc ...

For each models:

✓ imbalanced data

✓ Multi-class

✓ Interpretability

Prob. class labels

C.I on regression

Pros & cons

Outliers

obj  
- recap  
- applied  
- Interviews

Pre-processing (standardization)

Post-processing (calibration)

Good feature encoding

bias-variance tradeoff

Latency @ run-time

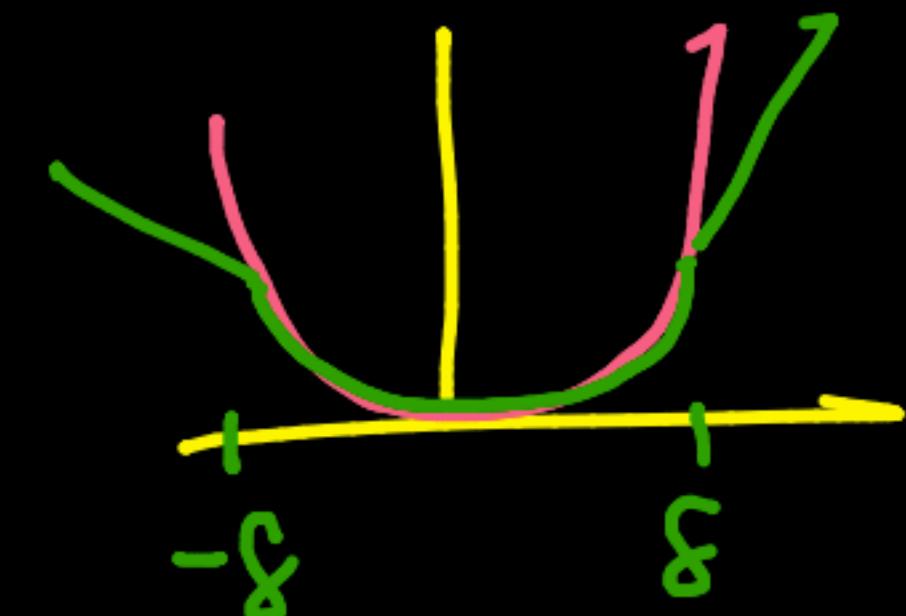
loss fn ...

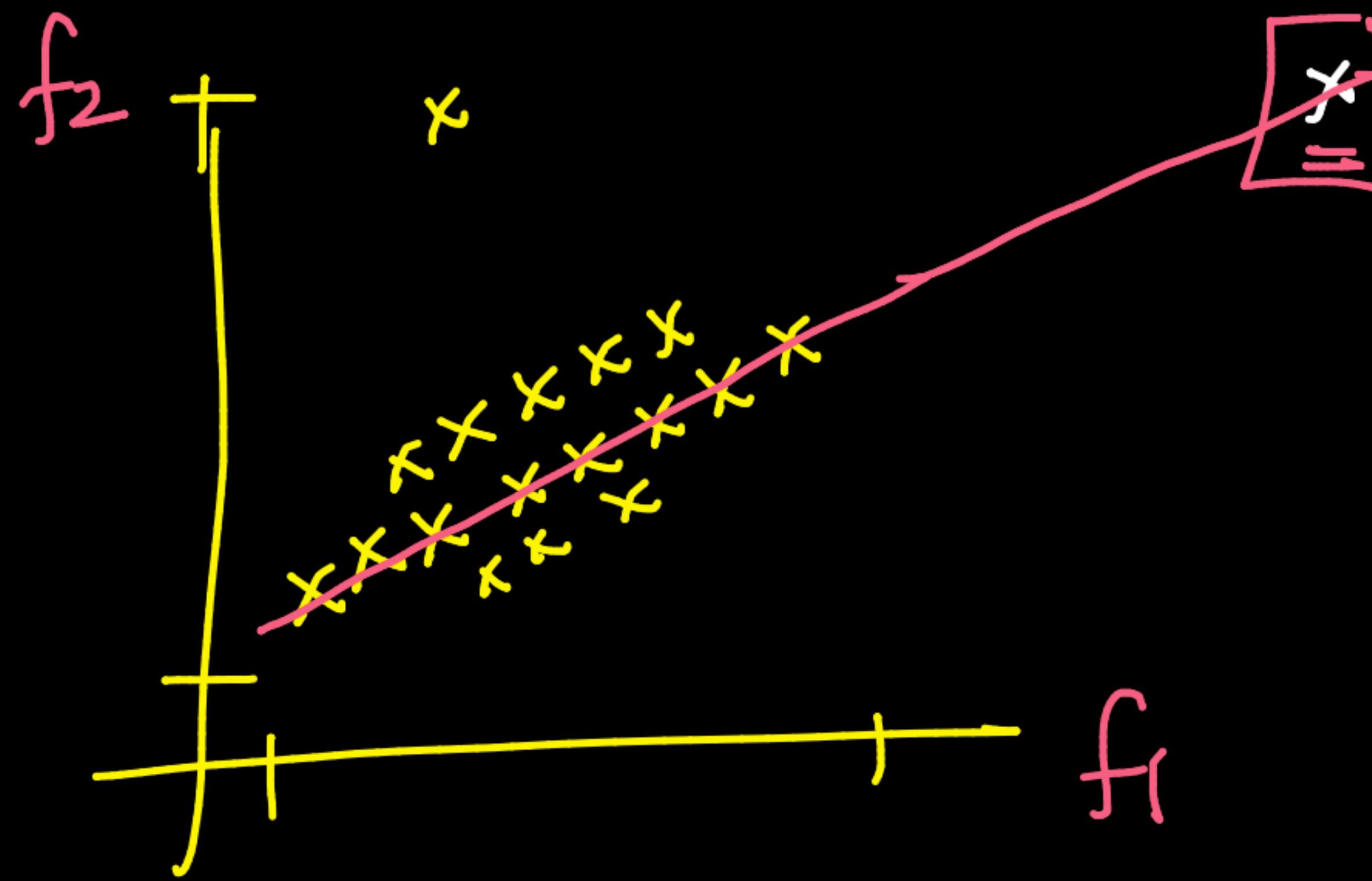
Data-drift ...

Linear-reg:

Loss: SQ-loss (MSE)  $\rightarrow$  MLE:  $y_i \sim N(\mathbf{w}^T \mathbf{x}_i + b, \epsilon_i)$ ,  
sq loss +  $\lambda \cdot \text{reg}$   $\rightarrow$  MAP: function of  $\sigma^2$   
bias-var: hyper-param  $\lambda$ : underfit  
 $\lambda \downarrow$ : overfit

Outliers: impacted  
fixes: RANSAC, MAE, Huber  
 $y_i - \hat{y}_i$





decision surface:  $\Pi^d$

↳ quadratic & higher order poly  $\rightarrow$  Curve or Surface in  $d$ -dim

$\hat{y}_i$  c.i. on  $y_i$

stat-assumptions

$$\checkmark \left\{ \begin{array}{l} \text{run} \\ \hat{w}^\top \hat{x}_q + b \end{array} \right. \pm 1.96 \hat{\sigma}$$

bootstrapping

disb of errors

$$\hat{y}_q \sim N(\hat{w}^\top \hat{x}_q + b, \hat{\sigma}^2)$$

$$\hat{\sigma}^2 = \text{MSE}_{\text{train}}$$

disadv

Train time

$M_1, M_2, \dots, M_{K^{100}}$   
eval time is large

$$N(0, \sigma^2)$$

Interpretability

$\rightarrow [w_j]$  on standardized data  
 $+ \text{ vs } -$

Collinear features  
 Correlated features

(later)

LIME, SHAP (DL)

adj R<sup>2</sup>

$$\log(f_i) = f_i \begin{vmatrix} f_i \\ f_j \\ f_{ij} \end{vmatrix}$$

↳ Reg  $\rightarrow$  Zerout

expensive

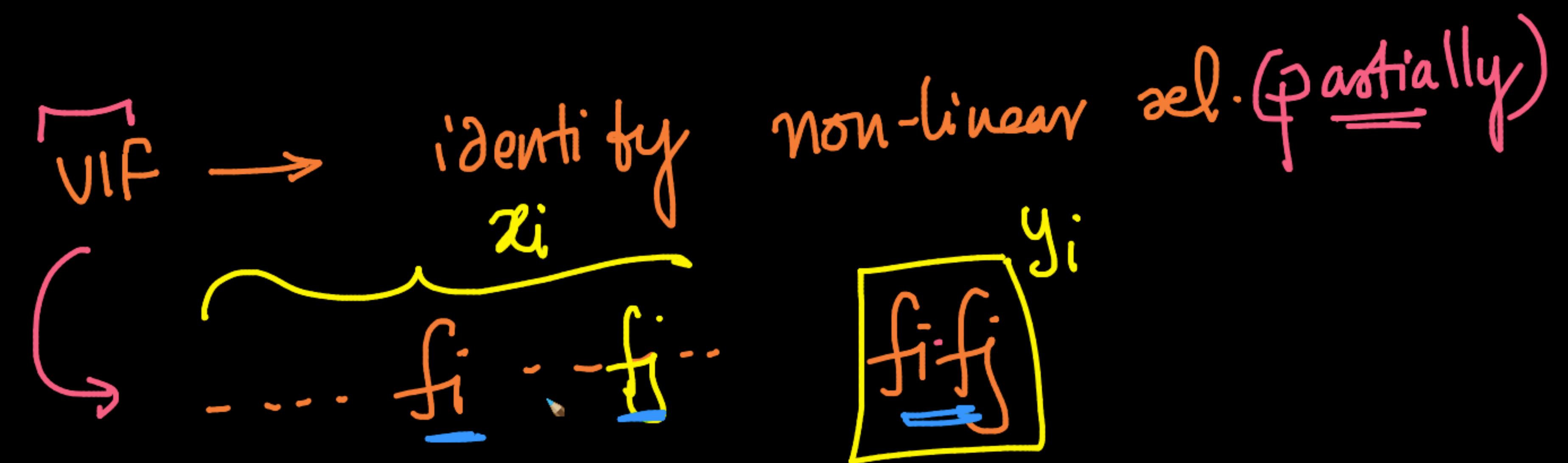
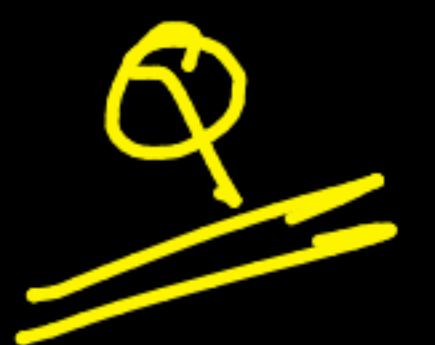
[Backward selection]

removing - each feature  
 $\Delta \text{Val loss}$

Pre-processing

↳ std · dala → interpretability of  $w_j$ 's  
→ optimization (SGD)

Post-processing:  $\hat{y}_i$ ; check  
 $f_i \sim Normal$  or not  
↳ c.l which can be trusted ...



Good feature encoding

{ cat-feat  $\rightarrow$  target;  
  Nom-feat  $\rightarrow$  std.

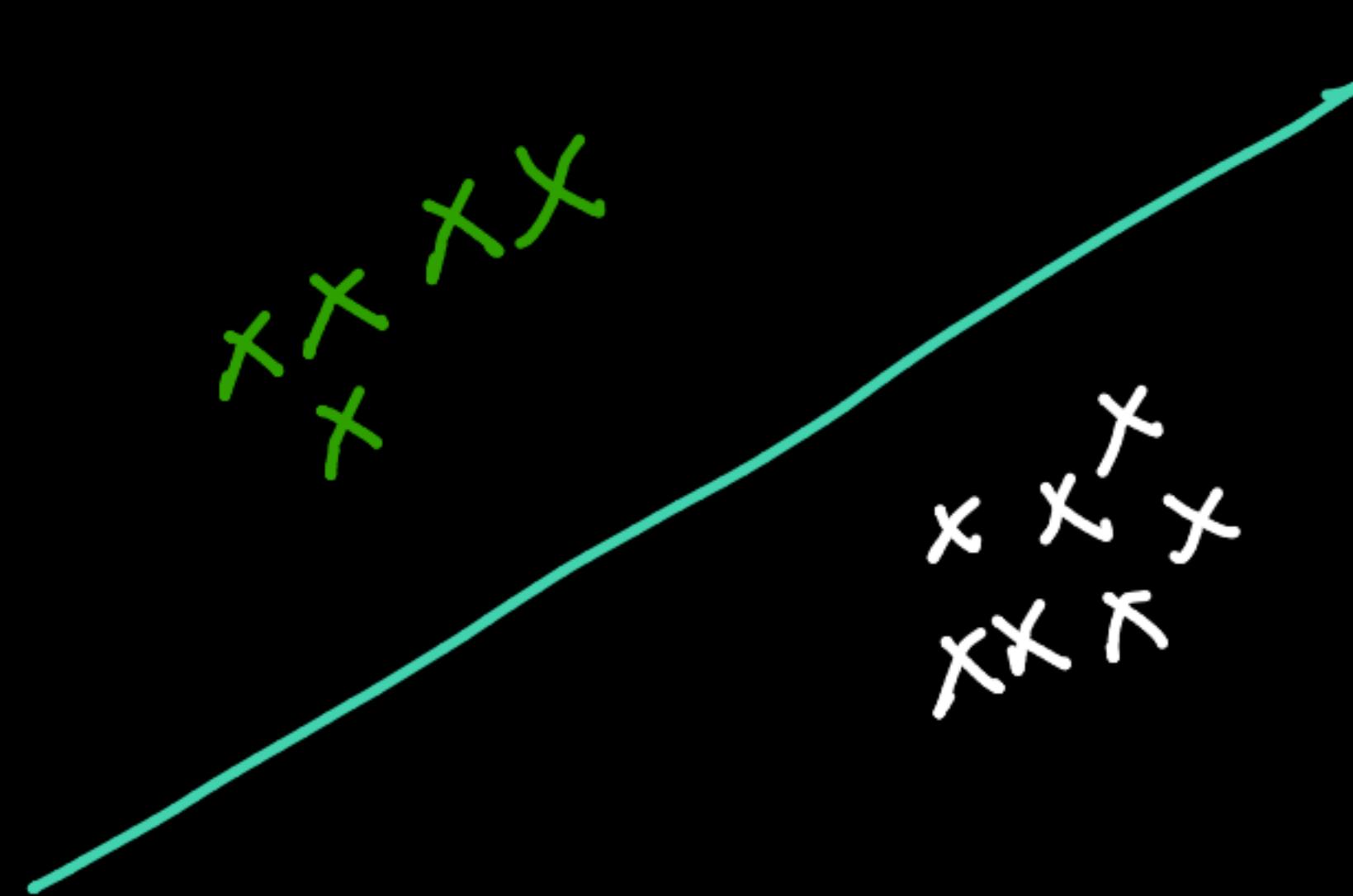
OHE  $\rightarrow$  dim increasing  
More "zoom" or  
flexibility to fit  
a  $\Pi_d$

$f_i$ : cat-feature

$c_1 c_2 \dots c_{20}$  (Categories)

L1 reg  
 $\rightarrow$  avoid  
overfitting

Curse of  
dim



Kernel-trick

Logistické reg:

Imbalanced →  
data

Logistické reg:

Imbalanced  
data

Upsampling;  
SMOTEj  
Down  
Sampling; Class-weights  
(hyper-param)

Multi-class  
Classfn → OVR; softmax  
(DL...)

Loss-fn: → log-loss(BCE)

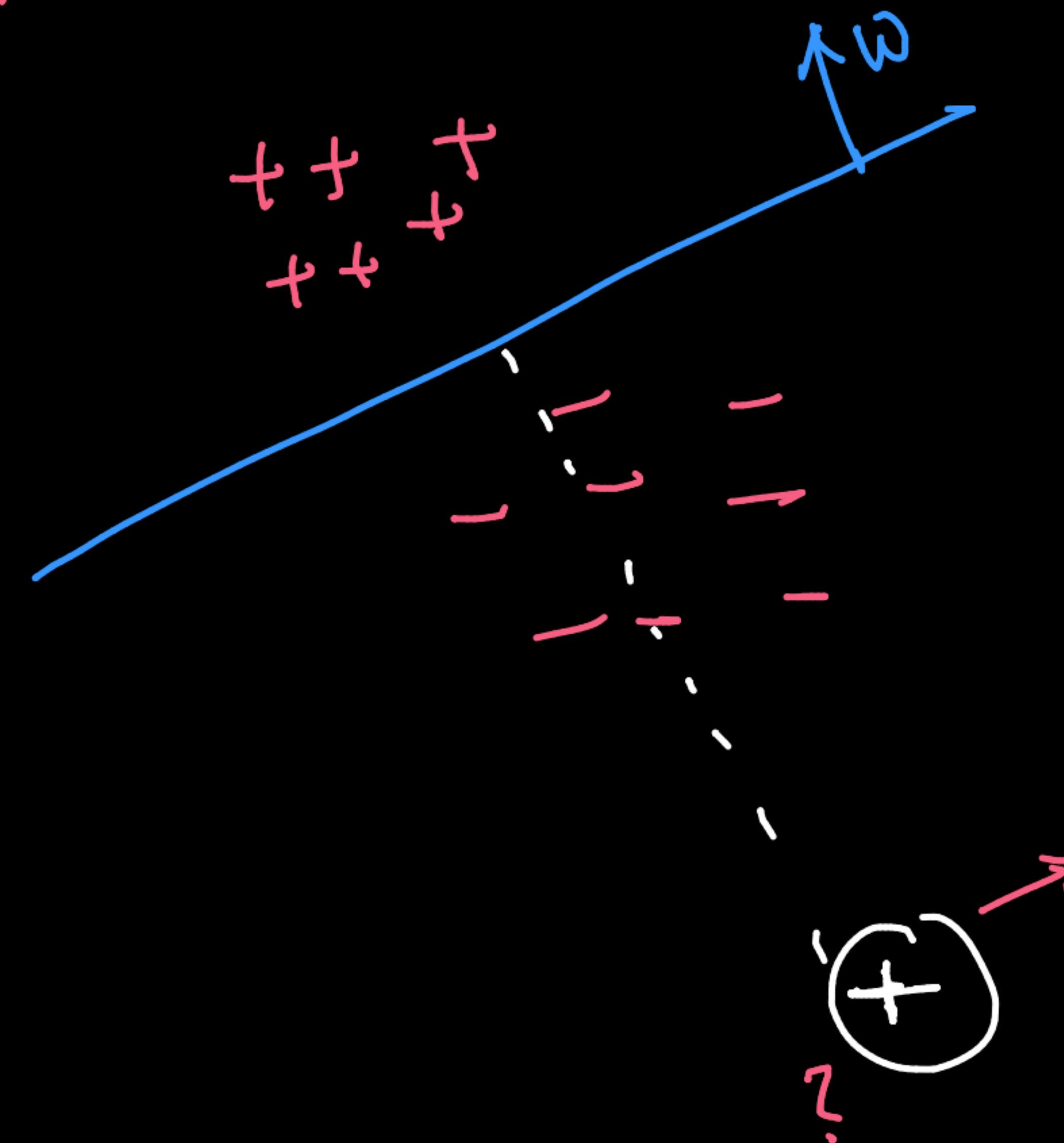
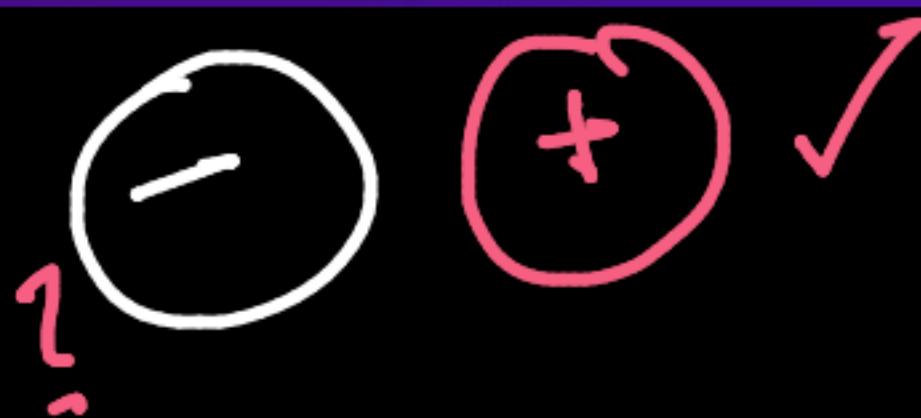
Interpretability  $\rightarrow$  same as lr-reg

feat. encoding  $\rightarrow$  "

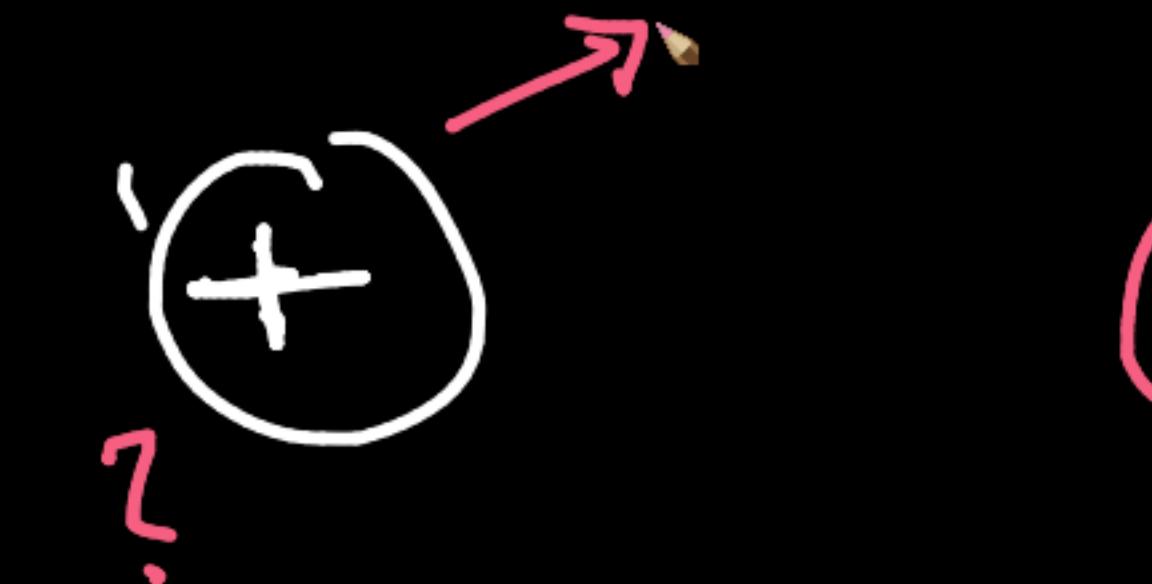
bias-var

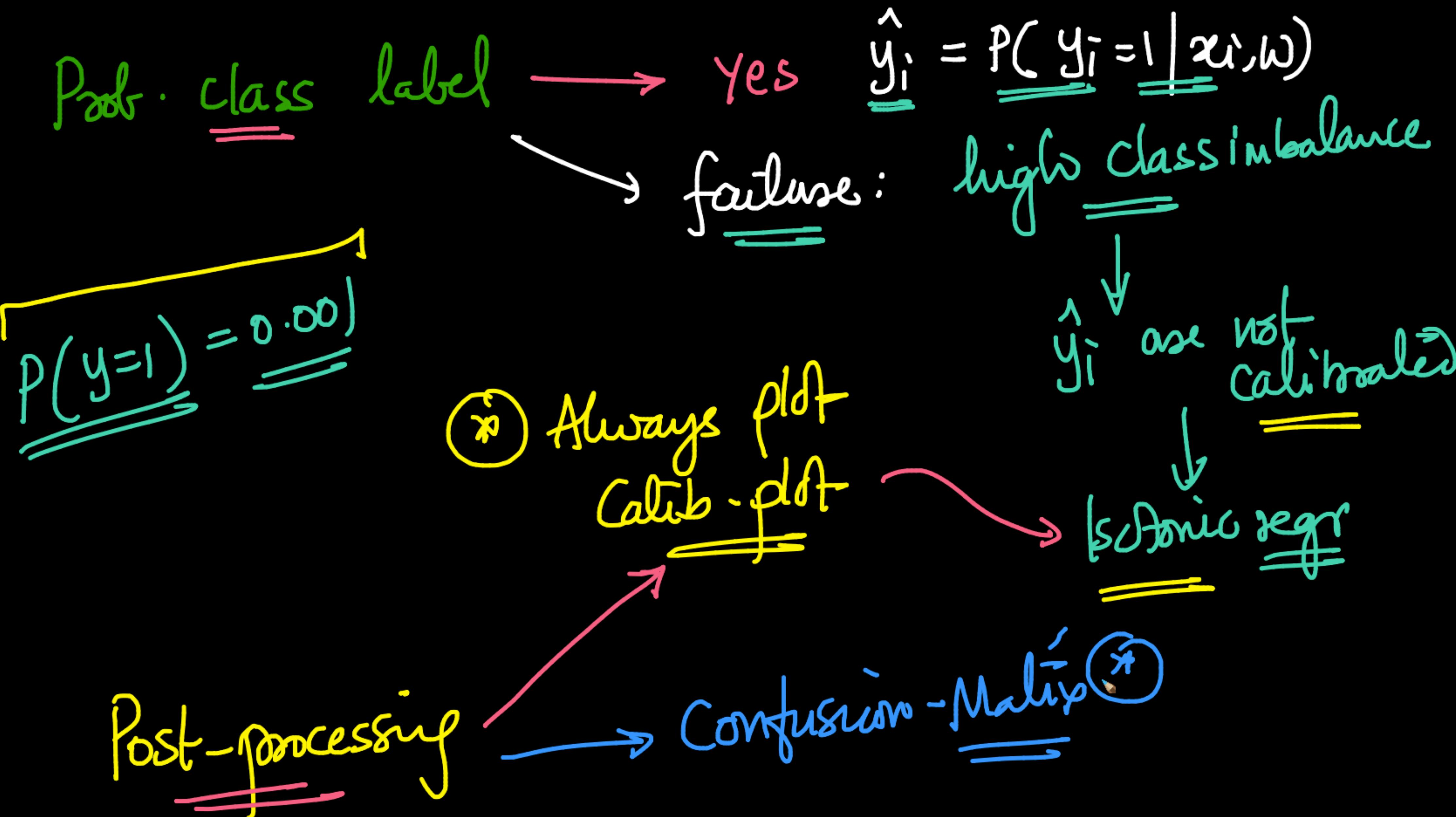
Outliers

" " "  
" " "  
RANSAC; change-loss for (research-)  
{ log-loss:  $y_i \log p_i + (1-y_i) \log (1-p_i)$   
                     $\hat{y}_i$



~~Sigmoid~~





Test-time:  $\mathcal{O}(d')$

$\equiv$

↑ # non-zero features

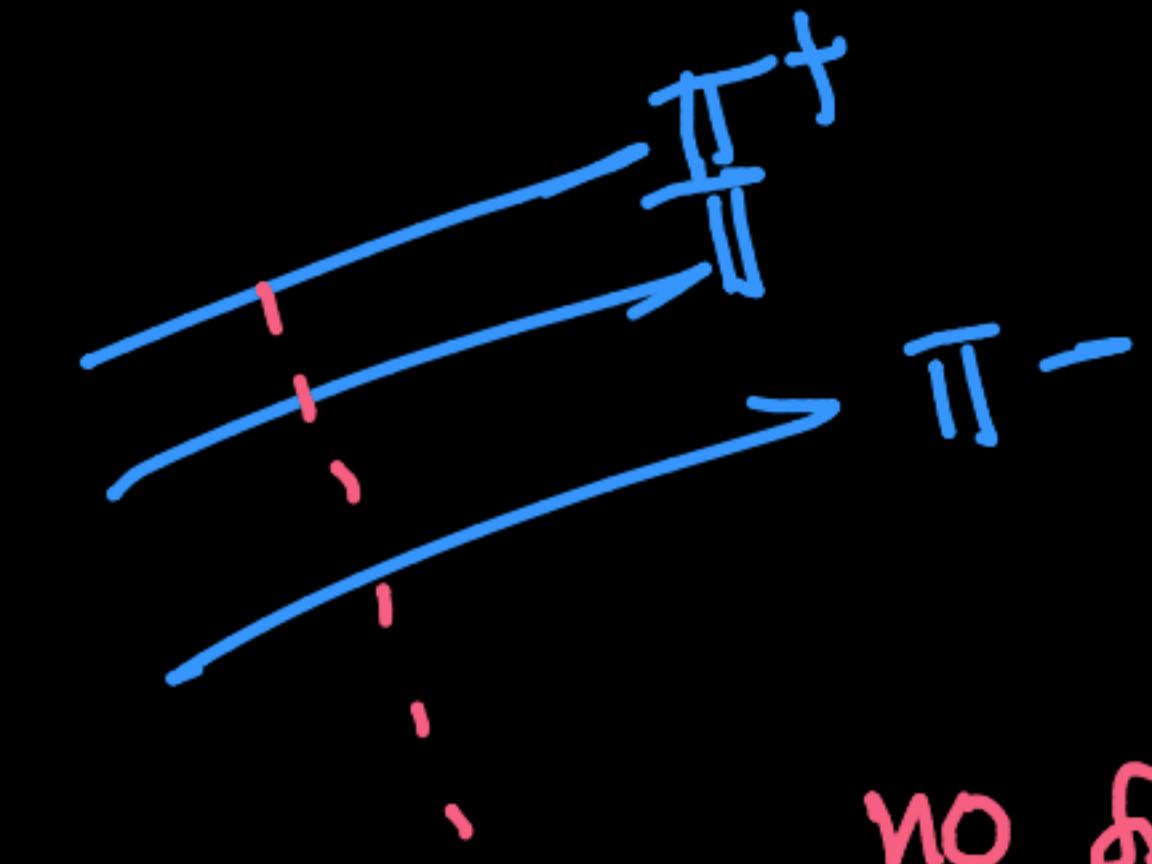
⊕ cons of lr-models

cannot automatically discover  
non-linear relationships

- L<sub>y</sub>-SVM
- Very similar to linear Models
- $\left[ \text{hinge-loss} + \frac{\|w\|_2^2}{2} \right]$  Primal
- Similar perf as linear Models
- Interpretability → ~~w<sub>j</sub>'s~~ (Scaling in Primal)
- ④ → Prob. class. labels  $\hookrightarrow P(y_i=1 | w, x_i)$

→ feat·enc  
↳ same as linear-model

→ outliers → " " "



no sigmoid  
to dampen fit



multi-class  $\rightarrow$  OVR

bias-var  $\rightarrow$  "C" =  $1/\lambda$

~~Suggestion~~



- think of all the cases - - - }
- WhatsApp - - - (weekend)



Minimizing (2) can be rewritten as a constrained optimization problem with a differentiable objective function in the following way.

For each  $i \in \{1, \dots, n\}$  we introduce a variable  $\zeta_i = \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i - b))$ . Note that  $\zeta_i$  is the smallest nonnegative number satisfying  $y_i(\mathbf{w}^\top \mathbf{x}_i - b) \geq 1 - \zeta_i$ .

Thus we can rewrite the optimization problem as follows

$$\text{minimize} \frac{1}{n} \sum_{i=1}^n \zeta_i + \lambda \|\mathbf{w}\|^2$$

subject to  $y_i (\mathbf{w}^\top \mathbf{x}_i - b) \geq 1 - \zeta_i$  and  $\zeta_i \geq 0$ , for all

This is called the *primal* problem.

## Dual [edit]

By solving for the Lagrangian dual of the above problem, one obtains the simplified problem

$$\text{maximize } f(c_1 \dots c_n) = \sum_{i=1}^n c_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i c_i (\mathbf{x}_i^\top \mathbf{x}_j) y_j$$

subject to  $\sum_{i=1}^n c_i y_i = 0$ , and  $0 \leq c_i \leq \frac{1}{2n\lambda}$  for all  $i$

This is called the *dual* problem. Since the dual maximization problem is a quadratic function of the  $c_i$  subject to linear constraints, it is efficiently solvable by [quadratic programming algorithms](#).

Here, the variables  $c_i$  are defined such that

$$\mathbf{w} = \sum_{i=1}^n c_i y_i \mathbf{x}_i.$$

Moreover,  $c_i = 0$  exactly when  $\mathbf{w}^\top \mathbf{x}_i = b$ . This means that  $\mathbf{w}$  is on the margin's boundary. It follows that  $\mathbf{w}$  can be written as a linear combination of the margin vectors.









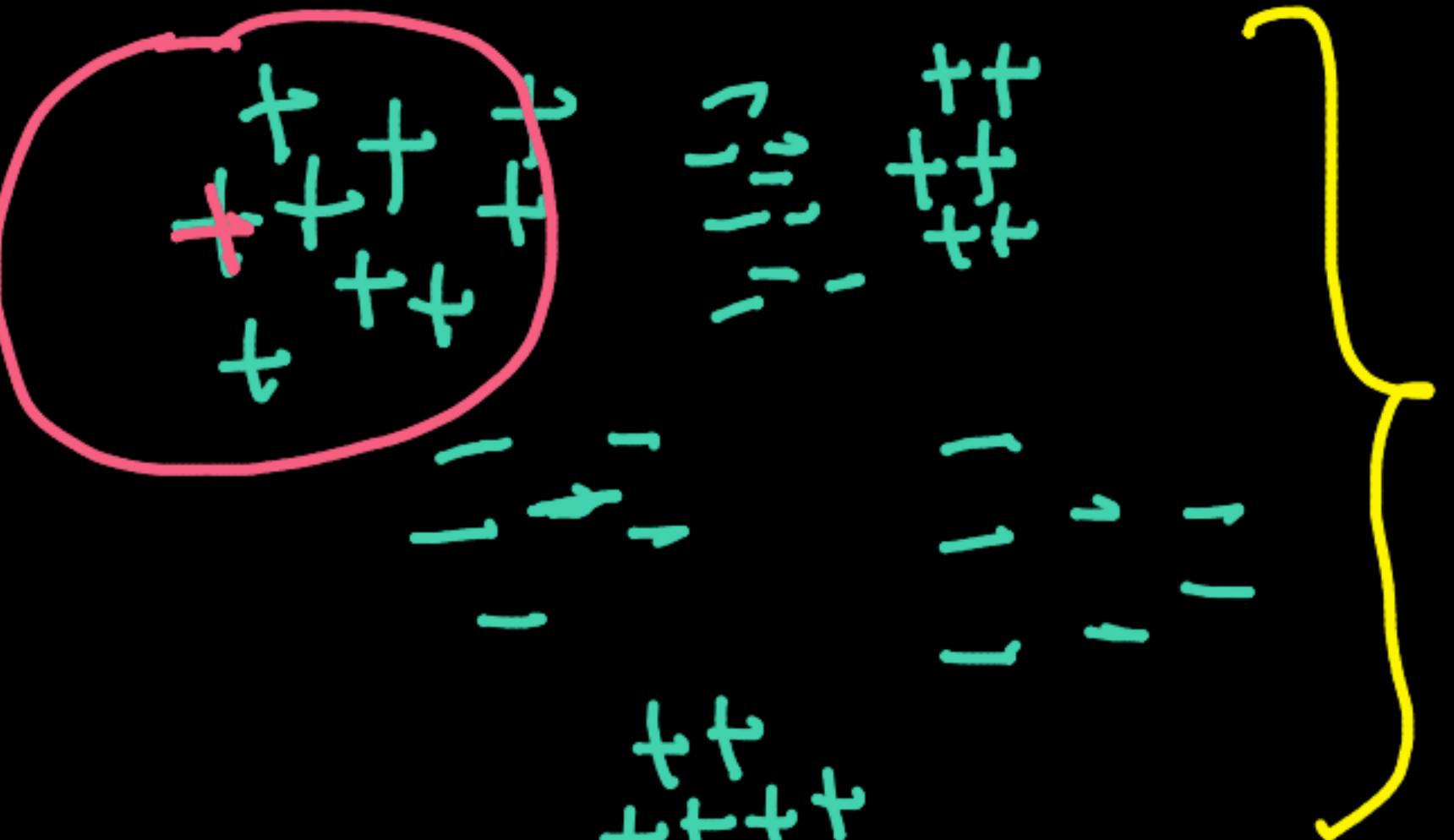
# Kernel-SVM

$\gamma$  is a DV:



→ no Interpolability ✓  
→ more time-complex... ✓  
→ #SUS ✓

→ Custom - domain specific  
are needed for good perf...



# RBF-SVM

(r)  
 $Q \uparrow$   
 $k \uparrow$  in CNN

This also helps to understand the [chi-squared distribution](#). Indeed, the (non-central) chi-squared distribution associated to a random point in the interval  $[0, 1]$  is the same as the distribution of the length-squared of a random point in the  $d$ -cube. By the law of large numbers, this distribution concentrates itself in a narrow band around  $d$  times the standard deviation squared ( $\sigma^2$ ) of the original derivation. This illuminates the chi-squared distribution and also illustrates that most of the volume of the  $d$ -cube concentrates near the boundary of a sphere of radius  $\sigma\sqrt{d}$ .

A further development of this phenomenon is as follows. Any fixed distribution on the [real numbers](#) induces a product distribution on points in  $\mathbb{R}^d$ . For any fixed  $n$ , it turns out that the difference between the minimum and the maximum distance between a random reference point  $Q$  and a list of  $n$  random data points  $P_1, \dots, P_n$  become indiscernible compared to the minimum distance:<sup>[12]</sup>

$$\lim_{d \rightarrow \infty} E \left( \frac{\text{dist}_{\max}(d) - \text{dist}_{\min}(d)}{\text{dist}_{\min}(d)} \right) \rightarrow 0.$$

This is often cited as distance functions losing their usefulness (for the nearest-neighbor criterion in feature-comparison algorithms, for example) in high dimensions. However, recent research has shown this to only hold in the artificial scenario when the one-dimensional distributions  $\mathbb{R}$  are [independent and identically distributed](#).<sup>[13]</sup> When attributes are correlated, data can become easier and provide higher distance contrast and the [signal-to-noise ratio](#) was found to play an important role, thus [feature selection](#) should be used.<sup>[13]</sup>

More recently, it has been suggested that there may be a conceptual flaw in the argument that contrast-loss creates a curse in high dimensions. Machine learning can be understood as the problem of assigning instances to their respective generative process of origin, with class labels acting as symbolic representations of individual generative processes. The curse's derivation assumes all instances are independent, identical outcomes of a single high dimensional generative process. If there is only one generative process, there would exist only one (naturally occurring) class and machine learning would be conceptually ill-defined in both high and low dimensions. Thus, the traditional argument that contrast-loss creates a curse, may be fundamentally inappropriate. In addition, it has been shown that when the generative model is modified to accommodate multiple generative processes, contrast-loss can morph from a curse to a blessing, as it ensures that the nearest-neighbor of an instance is almost-surely its most closely related instance. From this perspective, contrast-loss makes high dimensional distances especially meaningful and not especially non-meaningful as is often argued.<sup>[14]</sup>

## Nearest neighbor search [edit]

The effect complicates [nearest neighbor search](#) in high dimensional space. It is not possible to quickly reject candidates by using the difference in one coordinate as a lower bound for a distance based on all the dimensions.<sup>[15][16]</sup>

However, it has recently been observed that the mere number of dimensions does not necessarily result in difficulties,<sup>[17]</sup> since *relevant* additional dimensions can also increase the contrast. In addition, irrelevant ("noise") dimensions, however, reduce the contrast in the manner described above. In [time series analysis](#), where the data are inherently high-dimensional, distance functions also work reliably as long as the

This also helps to understand the [chi-squared distribution](#). Indeed, the (non-central) chi-squared distribution associated to a random point in the interval  $[0, 1]$  is the same as the distribution of the length-squared of a random point in the  $d$ -cube. By the law of large numbers, this distribution concentrates itself in a narrow band around  $d$  times the standard deviation squared ( $\sigma^2$ ) of the original derivation. This illuminates the chi-squared distribution and also illustrates that most of the volume of the  $d$ -cube concentrates near the boundary of a sphere of radius  $\sigma\sqrt{d}$ .

A further development of this phenomenon is as follows. Any fixed distribution on the [real numbers](#) induces a product distribution on points in  $\mathbb{R}^d$ . For any fixed  $n$ , it turns out that the difference between the minimum and the maximum distance between a random reference point  $Q$  and a list of  $n$  random data points  $P_1, \dots, P_n$  become indiscernible compared to the minimum distance:<sup>[12]</sup>

$$\lim_{d \rightarrow \infty} E \left( \frac{\text{dist}_{\max}(d) - \text{dist}_{\min}(d)}{\text{dist}_{\min}(d)} \right) \rightarrow 0.$$

This is often cited as distance functions losing their usefulness (for the nearest-neighbor criterion in feature-comparison algorithms, for example) in high dimensions. However, recent research has shown this to only hold in the artificial scenario when the one-dimensional distributions  $\mathbb{R}$  are [independent and identically distributed](#).<sup>[13]</sup> When attributes are correlated, data can become easier and provide higher distance contrast and the [signal-to-noise ratio](#) was found to play an important role, thus [feature selection](#) should be used.<sup>[13]</sup>

More recently, it has been suggested that there may be a conceptual flaw in the argument that contrast-loss creates a curse in high dimensions. Machine learning can be understood as the problem of assigning instances to their respective generative process of origin, with class labels acting as symbolic representations of individual generative processes. The curse's derivation assumes all instances are independent, identical outcomes of a single high dimensional generative process. If there is only one generative process, there would exist only one (naturally occurring) class and machine learning would be conceptually ill-defined in both high and low dimensions. Thus, the traditional argument that contrast-loss creates a curse, may be fundamentally inappropriate. In addition, it has been shown that when the generative model is modified to accommodate multiple generative processes, contrast-loss can morph from a curse to a blessing, as it ensures that the nearest-neighbor of an instance is almost-surely its most closely related instance. From this perspective, contrast-loss makes high dimensional distances especially meaningful and not especially non-meaningful as is often argued.<sup>[14]</sup>

## Nearest neighbor search [edit]

The effect complicates [nearest neighbor search](#) in high dimensional space. It is not possible to quickly reject candidates by using the difference in one coordinate as a lower bound for a distance based on all the dimensions.<sup>[15][16]</sup>

However, it has recently been observed that the mere number of dimensions does not necessarily result in difficulties,<sup>[17]</sup> since *relevant* additional dimensions can also increase the contrast. In addition,



DT + RF + GBDT  
(Tree based)

① Pre-processing: No standardization  
needed ...

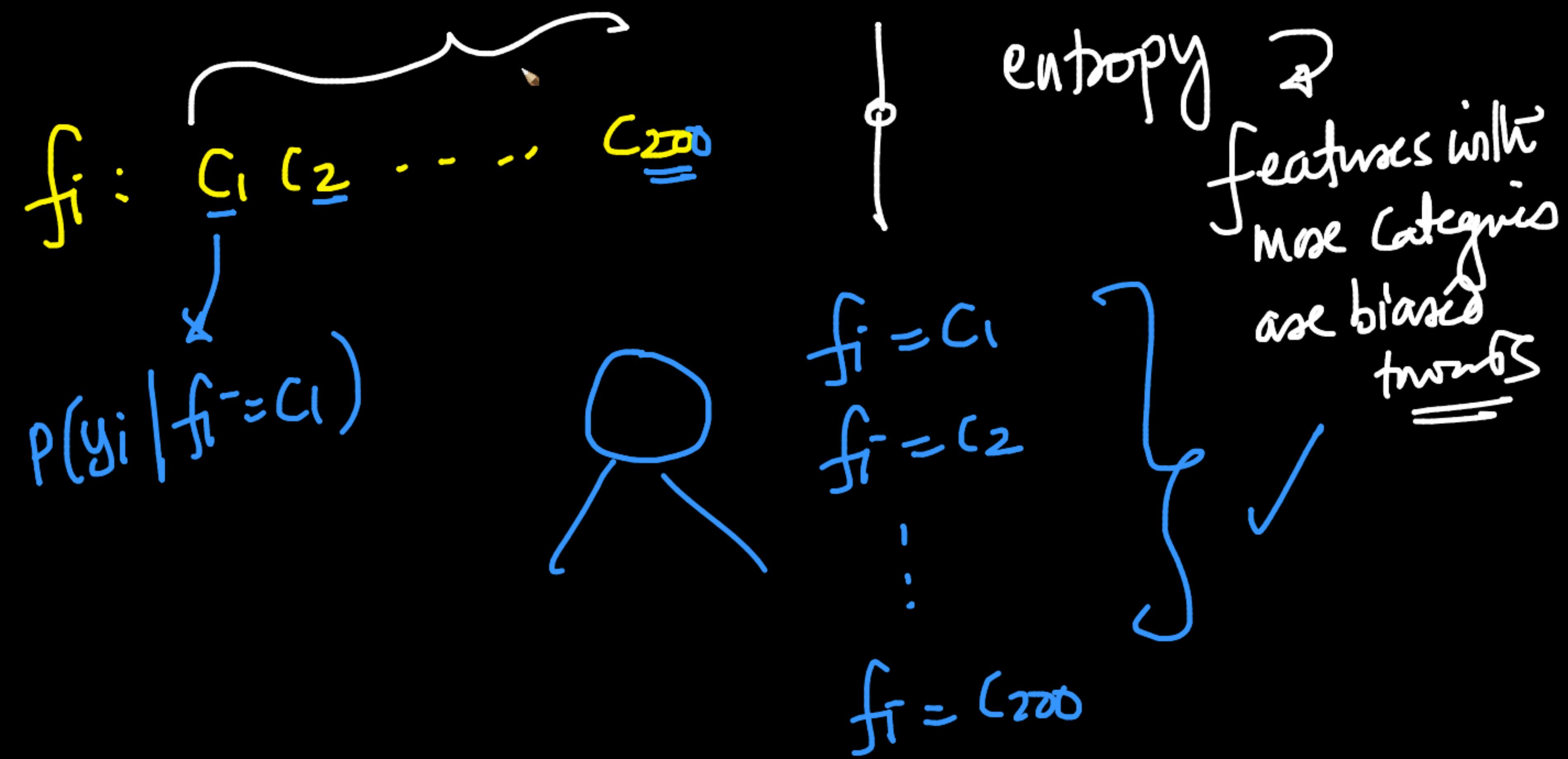
② Regression Using Trees:  $\rightarrow$  MSE



③ Feature-imp:  $\rightarrow$  weighted / normalized info.gain

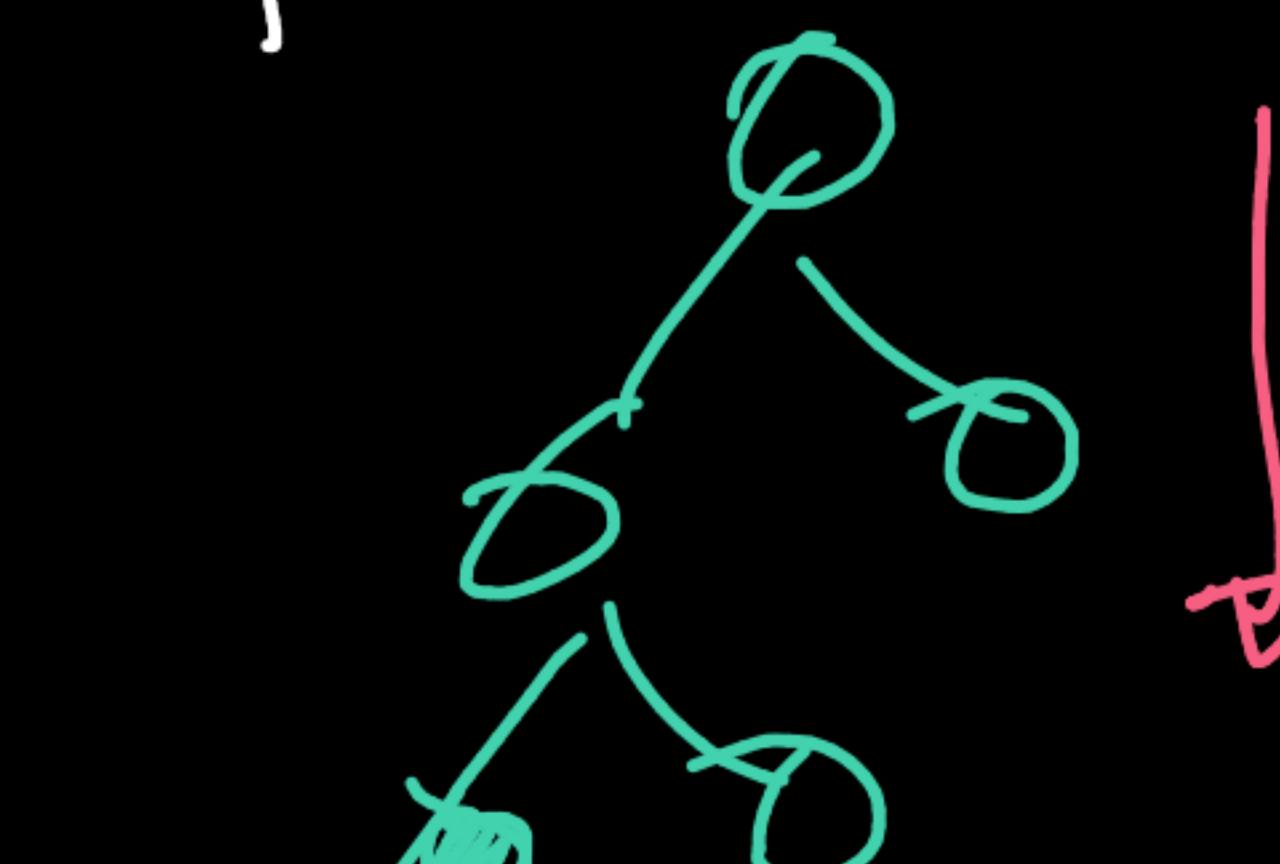
④ Outliers:  $\rightarrow$  not much

⑤ Feature-encoding:  $\rightarrow$  Cat-features  $\rightarrow$  Target ; other =  
 $\rightarrow$  Nom. feat  $\rightarrow$  as-is  
 $\rightarrow$  bins ...



C.I on  $\alpha_g \rightarrow$  bootstrapping  
Impurity in each leaf-node

Pos. class. labels  $\rightarrow$   
(toxic)



$$\{ (y_q=1) = \sum_L 2q \}$$
$$\begin{cases} 5: + \\ 1: - \end{cases}$$

# Naïve Bayes'-

- Laplace - Smoothing (bias-var)
- assumptions: conditional indep (polv.classes)

GBDT

→ any diff. loss

→ pseudo-code/internals

→ Xgboost vs LightGBM vs CatBoost  
vs GBDT

Tabular (default)

to find some function  $\hat{F}(x)$  that best approximates the output variable from the values of input variables. This is formalized by introducing some [loss function](#)  $L(y, F(x))$  and minimizing it in expectation:

$$\hat{F} = \arg \min_F \mathbb{E}_{x,y} [L(y, F(x))].$$

The gradient boosting method assumes a real-valued  $y$ . It seeks an approximation  $\hat{F}(x)$  in the form of a weighted sum of  $M$  functions  $h_m(x)$  from some class  $\mathcal{H}$ , called base (or [weak](#)) learners:

$$\hat{F}(x) = \sum_{m=1}^M \gamma_m h_m(x) + \text{const.}$$

We are usually given a training set  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  of known sample values of  $x$  and corresponding values of  $y$ . In accordance with the [empirical risk minimization](#) principle, the method tries to find an approximation  $\hat{F}(x)$  that minimizes the average value of the loss function on the training set, i.e., minimizes the empirical risk. It does so by starting with a model, consisting of a constant function  $F_0(x)$ , and incrementally expands it in a [greedy](#) fashion:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma),$$

$$F_m(x) = F_{m-1}(x) + \arg \min_{h_m \in \mathcal{H}} \left[ \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + h_m(x_i)) \right],$$

for  $m \geq 1$ , where  $h_m \in \mathcal{H}$  is a base learner function.

Unfortunately, choosing the best function  $h_m$  at each step for an arbitrary loss function  $L$  is a computationally infeasible optimization problem in general. Therefore, we restrict our approach to a simplified version of the problem.

The idea is to apply a [steepest descent](#) step to this minimization problem (functional gradient descent).

The basic idea behind the steepest descent is to find a local minimum of the loss function by iterating on  $F_{m-1}(x)$ . In fact, the local maximum-descent direction of the loss function is the negative gradient.<sup>[10]</sup>

Hence, moving a small amount  $\gamma$  such that the linear approximation remains valid:

$$F_m(x) = F_{m-1}(x) - \gamma \sum_{i=1}^n$$

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

In the discrete case however, i.e. when the set  $\mathcal{H}$  is finite, we choose the candidate function  $h$  closest to the gradient of  $L$  for which the coefficient  $\gamma$  may then be calculated with the aid of [line search](#) on the above equations. Note that this approach is a heuristic and therefore doesn't yield an exact solution to the given problem, but rather an approximation. In pseudocode, the generic gradient boosting method is:<sup>[5][2]</sup>

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

3. Compute multiplier  $\gamma_m$  by solving the following [one-dimensional optimization](#) problem:

$$\boxed{\gamma_m} = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

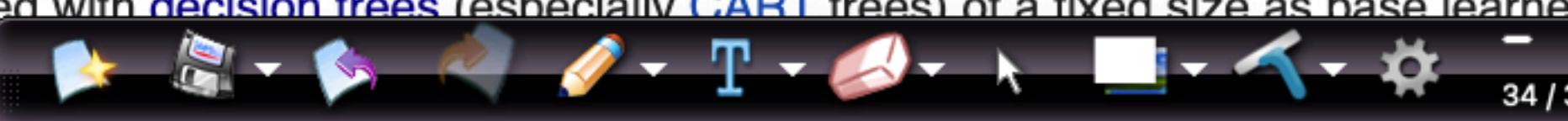
4. Update the model:

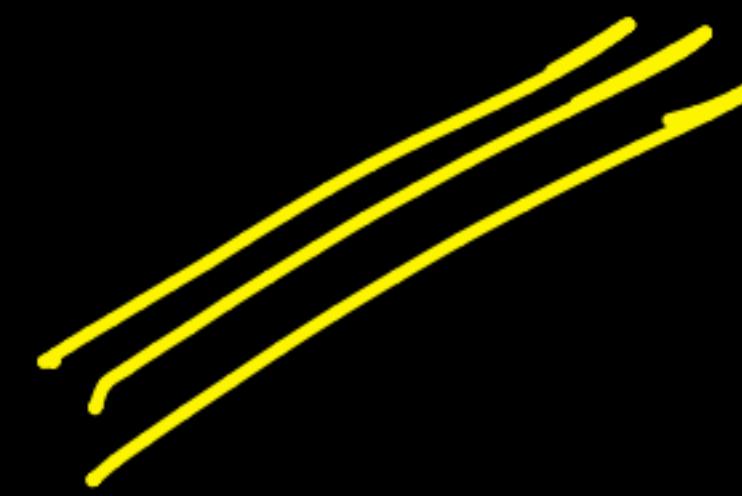
$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output  $F_M(x)$ .

## Gradient tree boosting [\[edit\]](#)

Gradient boosting is typically used with [decision trees](#) (especially [CART](#) trees) of a fixed size as base learners. For this special case, Friedman proposes a modification to gradient boosting method whi



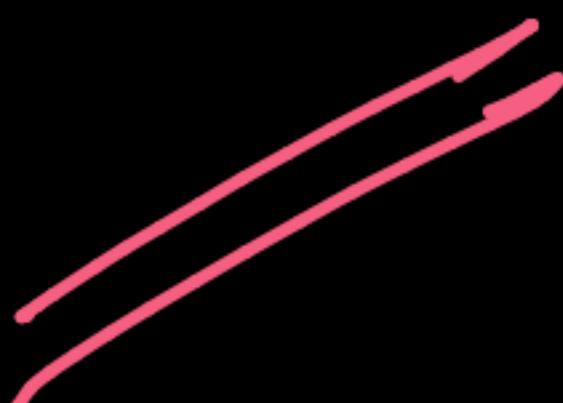


Time & Space - Compl

→ Run-time → ...  
→ Training time  
    ↓

Linear seg:

"i" ← learning-rate  
    ↑  
    ↓  
 $\frac{\text{GD}}{(n)} \mid \text{SGD} \mid \text{mini batches } = \underline{(k)}$



# Feat. Imp in Linear Reg

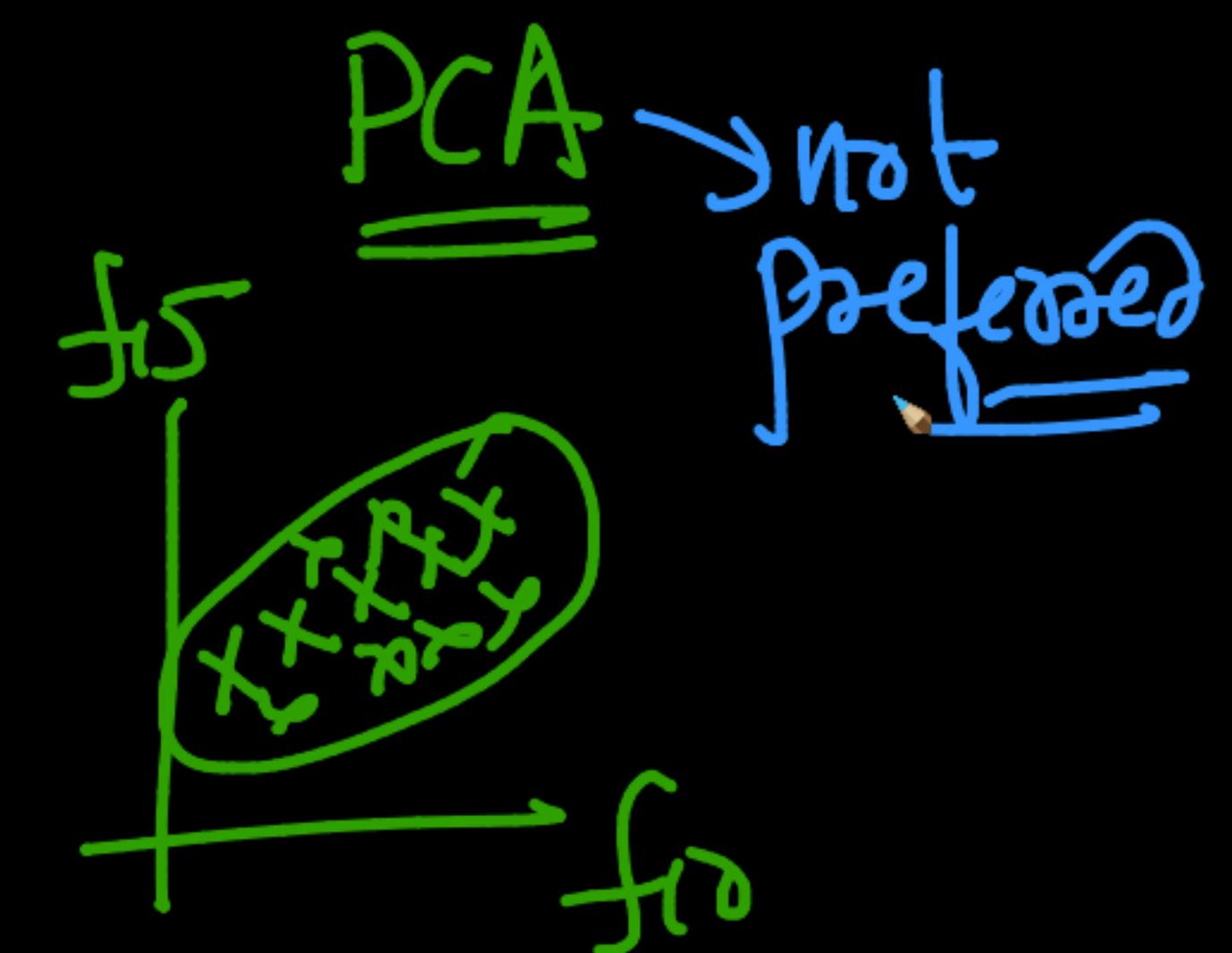
$d=50$

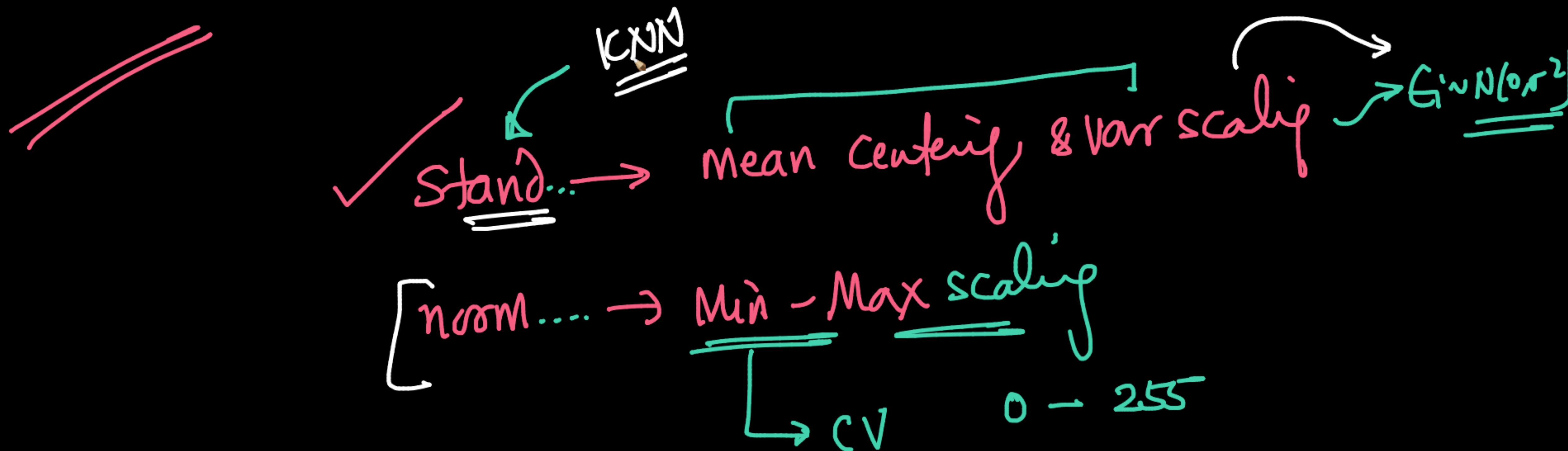
VIF can easily handle this

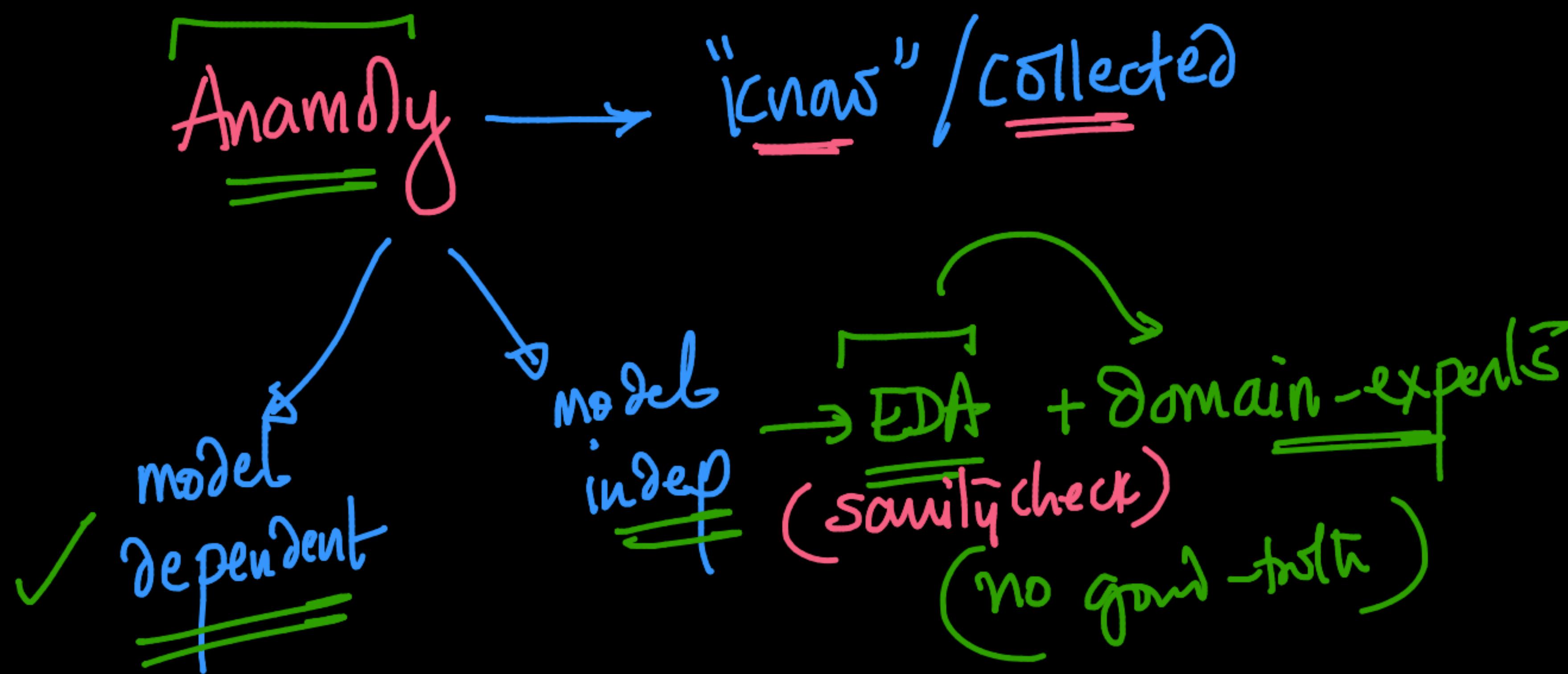
$\dots \underline{f_{10}} \dots \underline{f_{15}} \dots \underline{f_{50}}$

(Q1)  $d \xrightarrow{\text{PCA}} d'$

(Q2)







Anomaly / outliers / Novelties

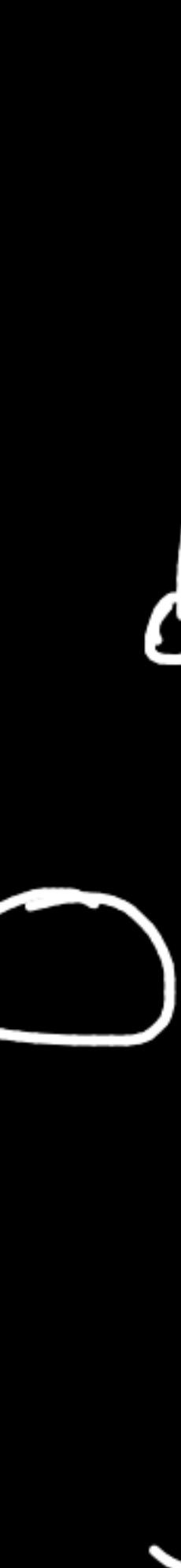


generative

Naive Bayes

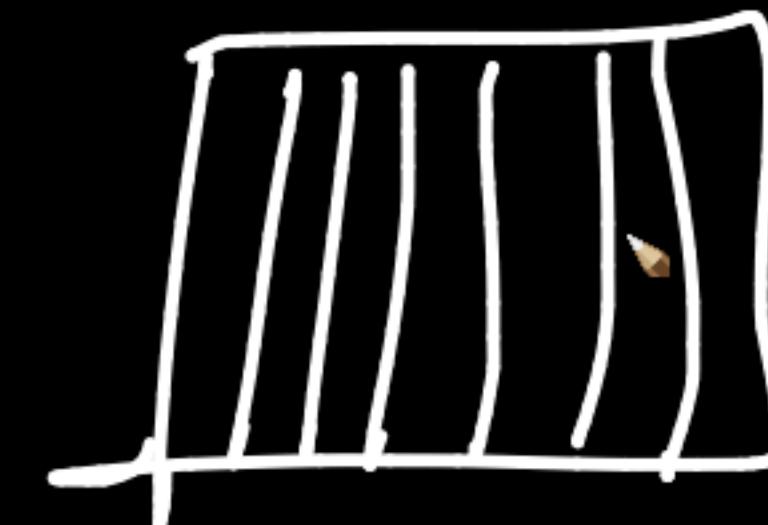
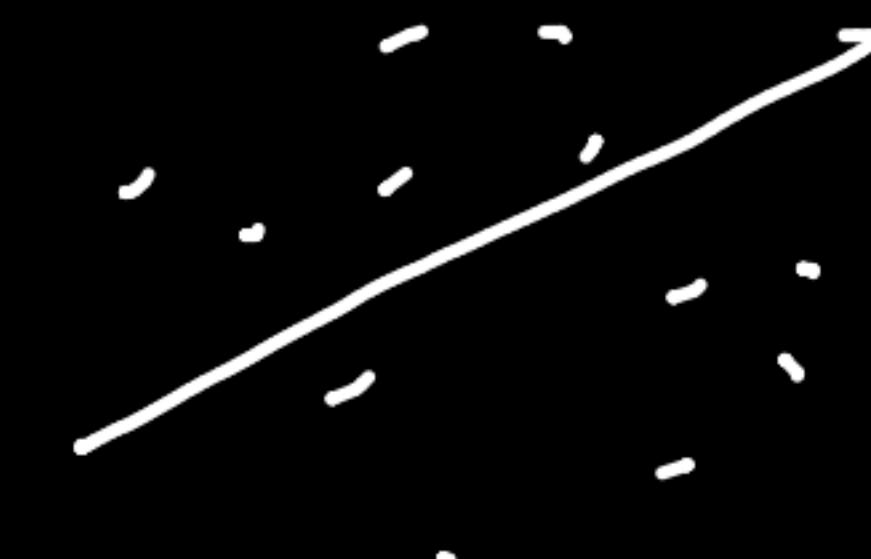
GMM

GAN (cv) / Dall-E

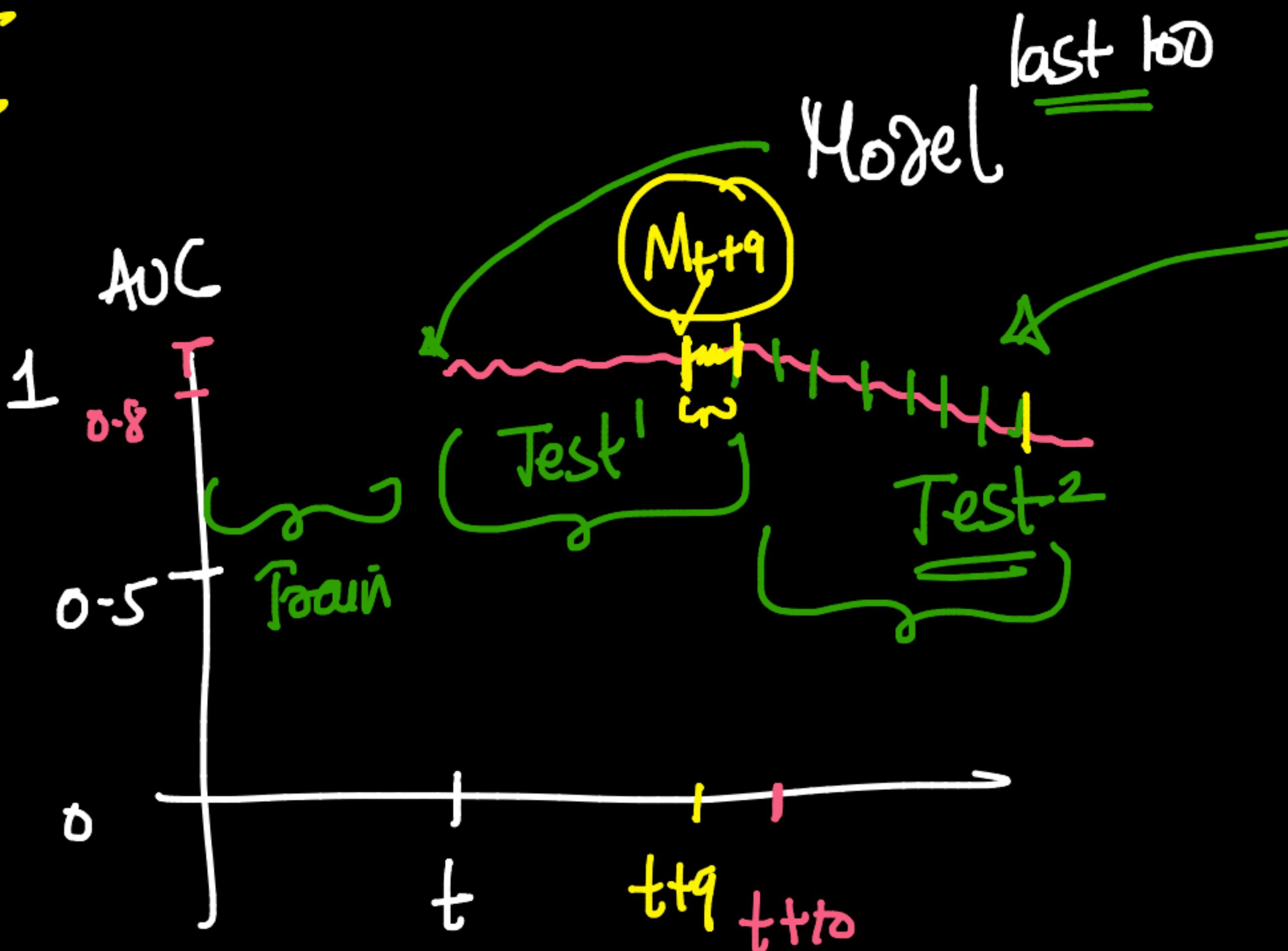


discriminative

DT e.g.



Data drift



$\mathcal{D}^1 \dots \mathcal{D}^{100} \rightarrow w=1$

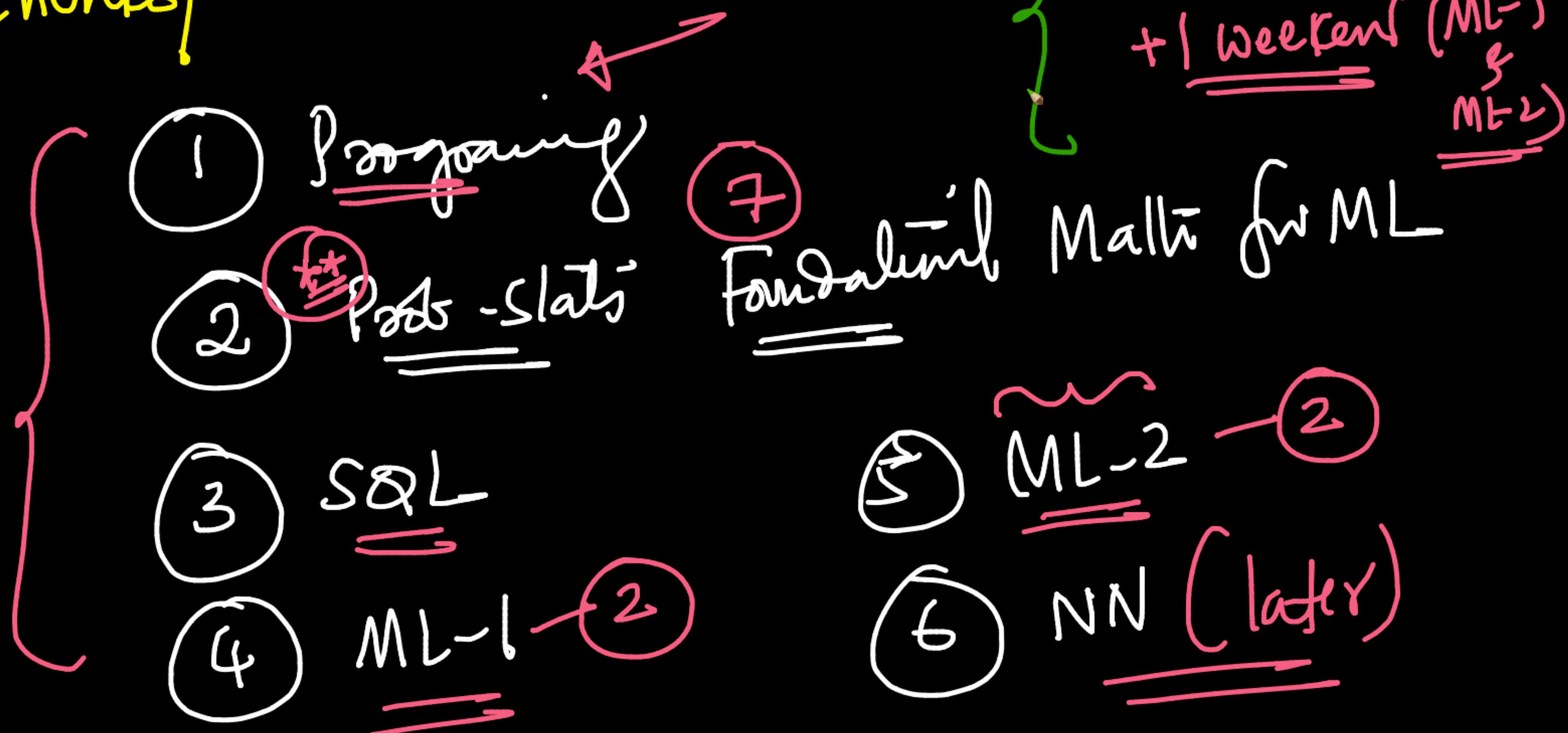
$M_{1-10} \rightarrow d_{100}$

$\rightarrow d_{101}$

{  
  de-Train  
  init:  $M_{1-10}$   
  Update:  $\mathcal{D}_{\text{Oversampled}}^{101}$

$\tilde{\mathcal{D}}^{101} \rightarrow \text{oversample}$   
 $\underline{\mathcal{D}^{101}} \rightarrow \text{SMOTE}$   
 $w=10$

# Chunks | Modules:



Curse of dimensionality - Wikipedia | Support-vector machine - Wikipedia | Gradient boosting - Wikipedia | Computational complexity of machine learning algorithms | 3blue1brown Fourier transforms

youtube.com/results?search\_query=3blue+1+brown+fourier+transforms

3blue1brown Fourier transforms

FILTERS

Did you mean: 3 blue 1 brown fourier transform

Signal 

Winding 

Transform 

But what is the Fourier Transform? A visual introduction.

7.8M views • 4 years ago

3Blue1Brown CC

If you want to check it out, I feel compelled to warn you that it's not the most well-documented tool, and has many other quirks you ...

What's that? | "Almost" Fourier transform? | Inverse Fourier?

3 moments

$n = 10$   $n = 50$   $n = 250$

Drawn with circles 

24:47

But what is a Fourier series? From heat flow to drawing with circles | DE4

14M views • 3 years ago

3Blue1Brown CC

Small correction: at 9:33, all the exponents should have a  $\pi^2$  in them. If you're looking for more Fourier Series content online, ...

Oooh, circle drawings | The heat equation | Interpreting the infinite sum | To the complex plane |...

7 chapters

Pure Fourier series animation montage

602K views • 3 years ago

3Blue1Brown CC

If you're curious about the number of vectors used for each animation: - Eighth note: 100 - Sigma: 200 - Britain: 500 - Fourier ...

Pure Fourier transforms

43 / 43

There's a key way in which the description I gave of the trade-off in Doppler radar differs from reality. Since the speed of light is so ...

MORE FROM YOUTUBE

Home Explore Shorts Subscriptions Library History Watch later Liked videos

SUBSCRIPTIONS

Music Sports Gaming Movies

EXPLORE

Movies Gaming Live Fashion & Beauty Learning Sports