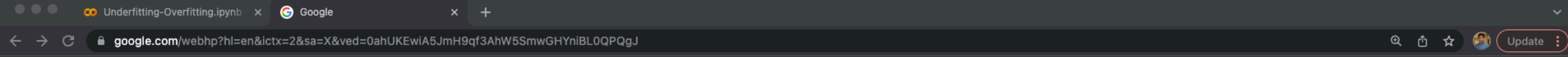


Underfitting - Overfitting

Bias- Variance Tradeoff - 2

Linear Regression

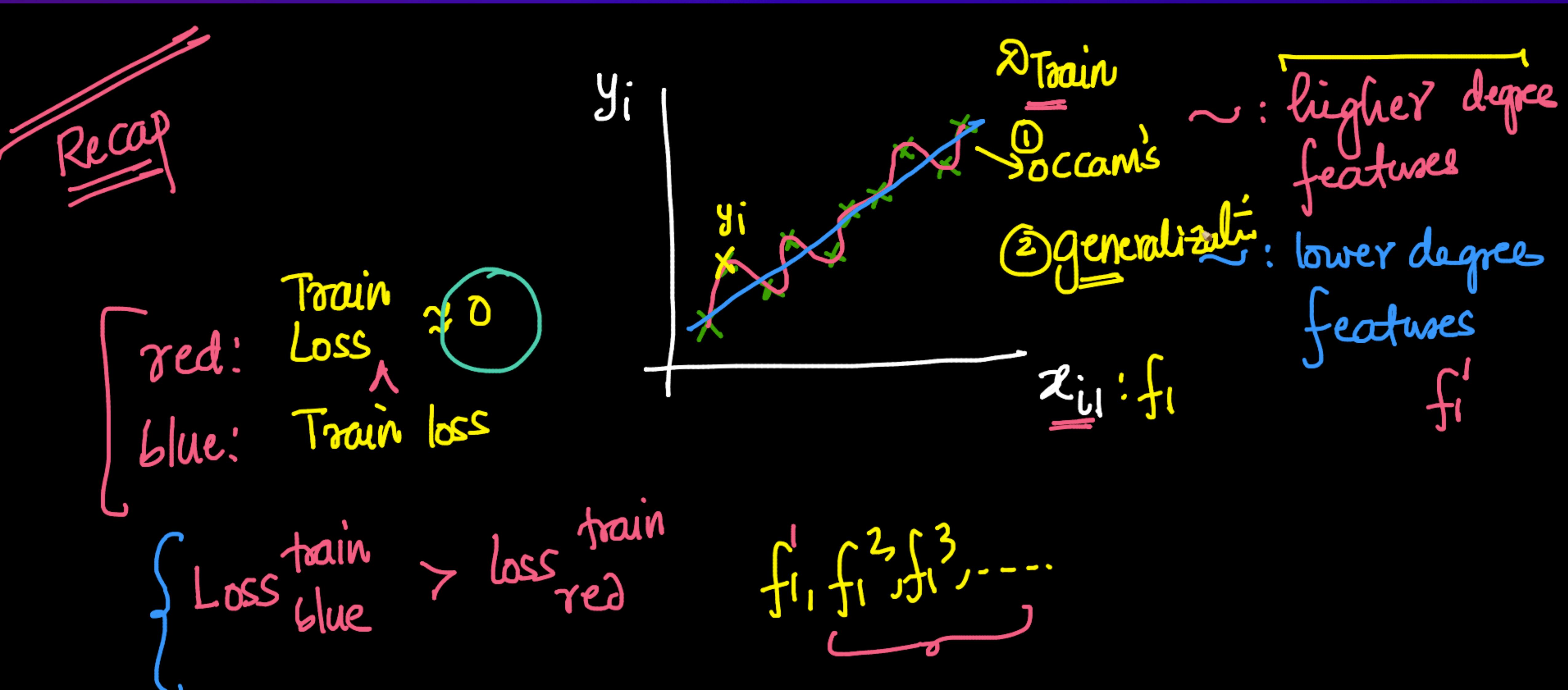
generic

[Gmail](#)[Images](#)

India

Agenda:

- Underfitting vs overfitting Tradeoff (examples)
- Intuition & Math (optimization)
- { - end-end Linear Regression
 - Hyper-param tuning
 - Code (simulation)



$$\text{Loss} = \sum_{i=1}^n l_i \rightarrow (y_i - \hat{y}_i)^2 \rightarrow \text{squared loss}$$

(linear-reg)

$$\min_{w_j} \sum_{i=1}^n [y_i - (w^T x_i + w_0)]^2$$

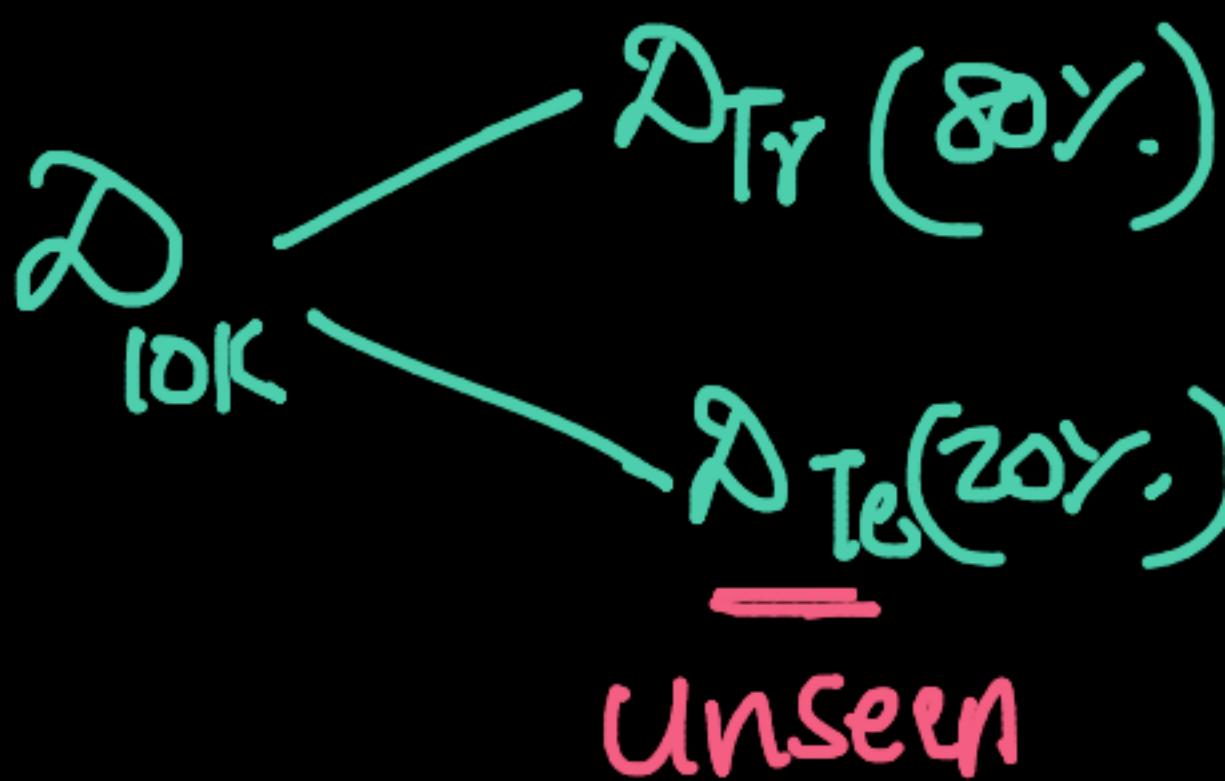
① Generalization: ML
Train - loss

$D_{TR} \rightarrow$ Model

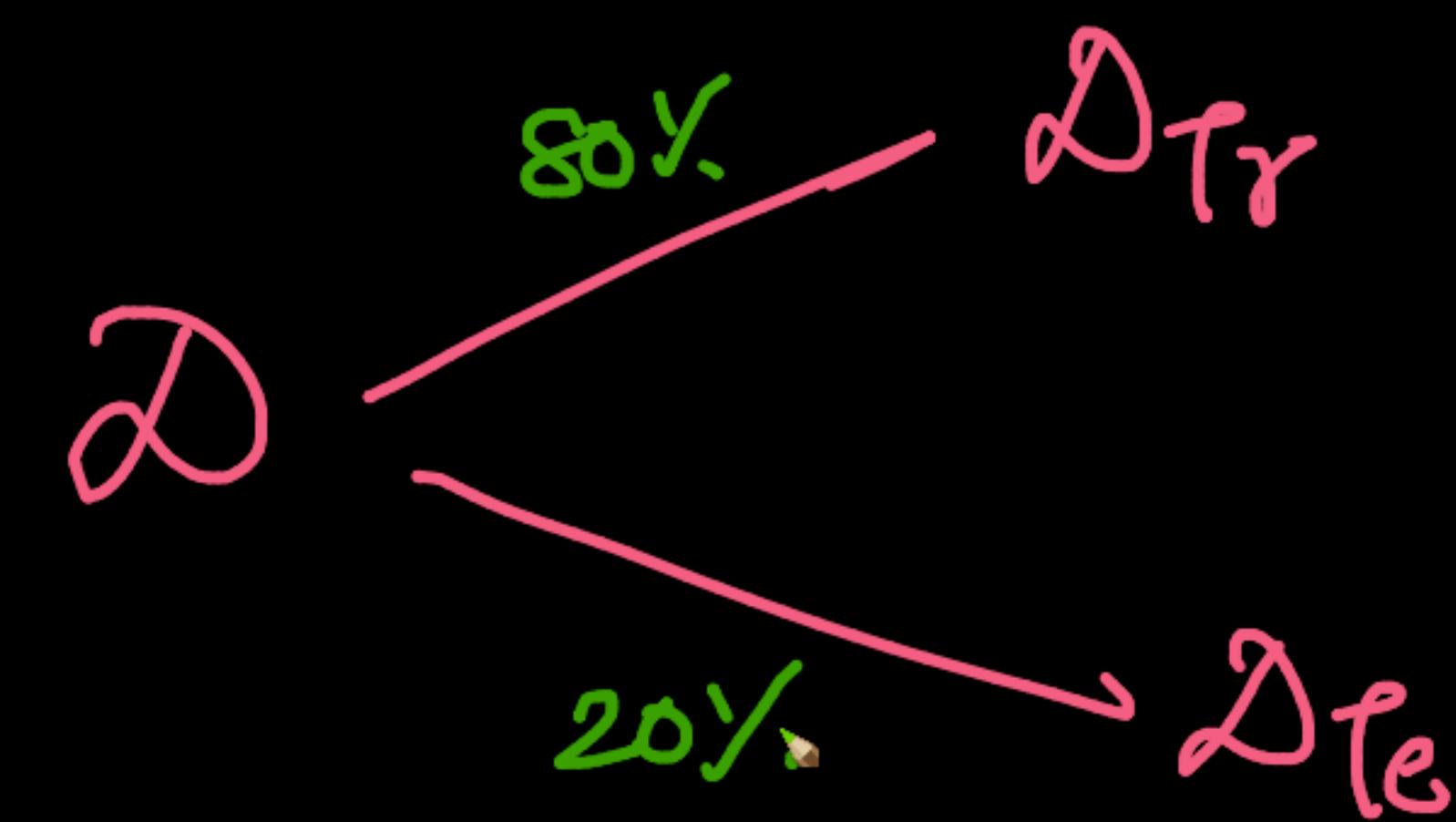
unseen { D_{Te} \rightarrow Model \rightarrow R-squared

② Occam's razor: science (ML)

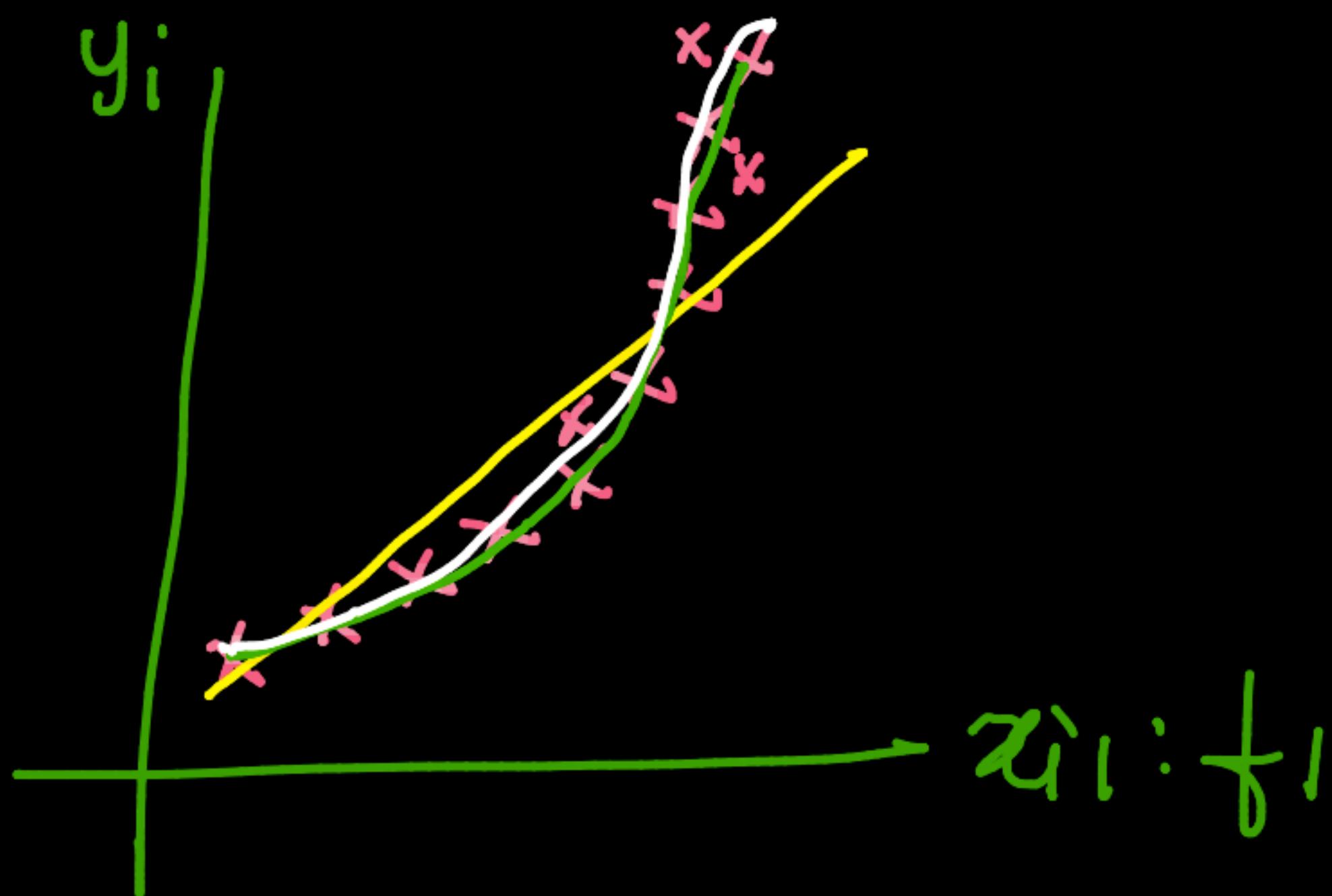
↳ Pick the simplest model that fits/explain your data well



Real-world: All data is messy



Q



- ① Linear : f_1^1
- ② Quadratic : f_1, f_1^2
- ③ Cubic : f_1, f_1^2, f_1^3

simplest

✓ Occam's razor

✓ [Gen:]

Linear

#

quad

L

complex
cubic

VL

4
bicubic

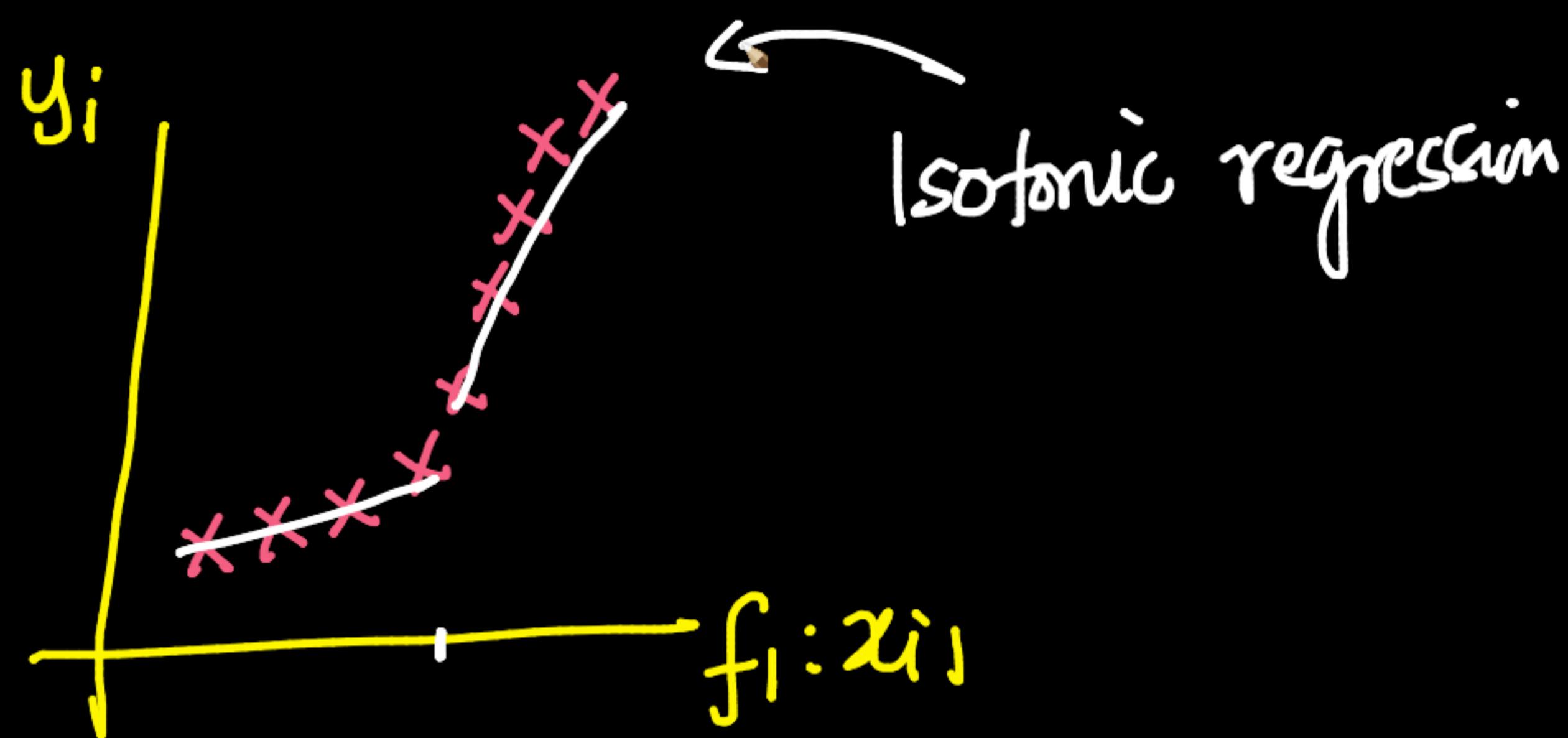
5th

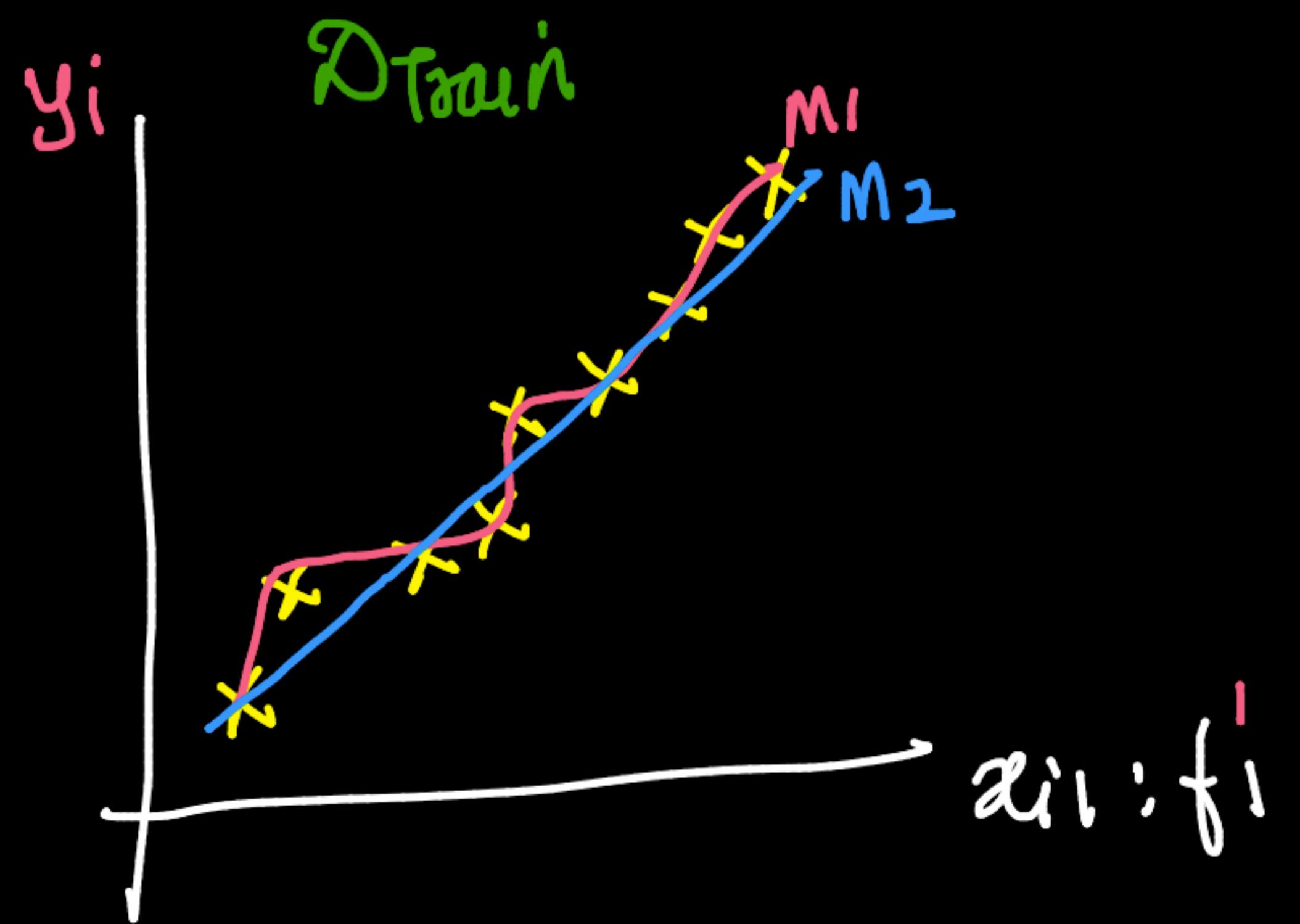
Linear-model: solving optimzn

(one-plane =)

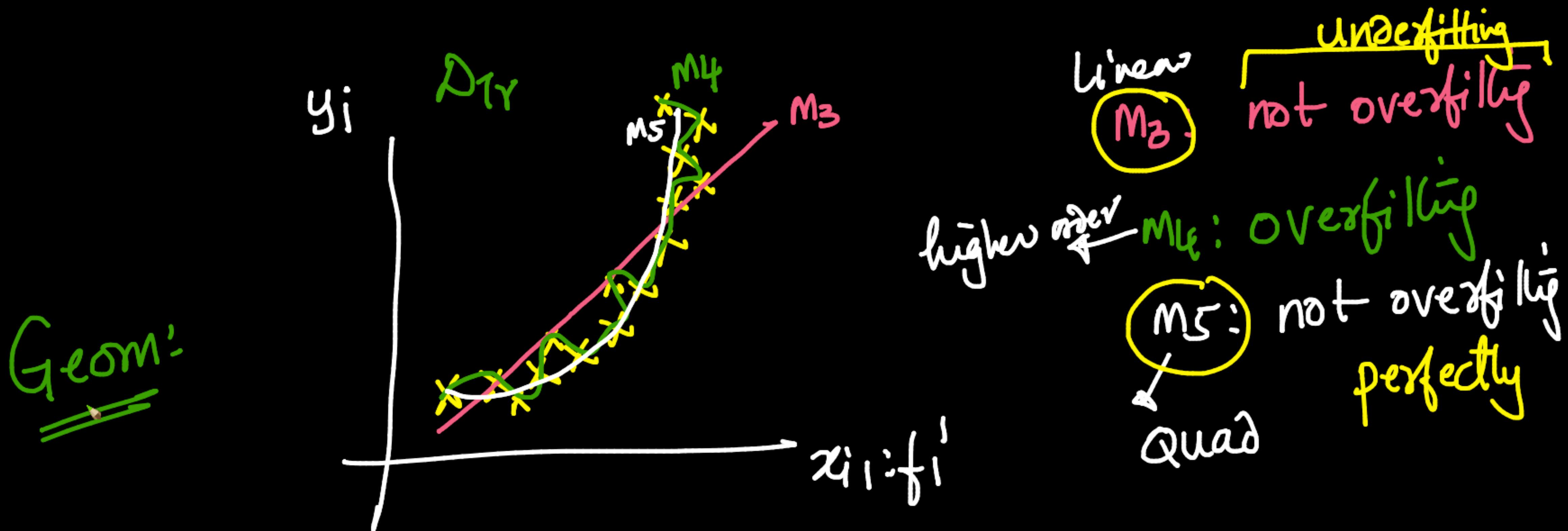
$$\min_{w_j} \sum_{i=1}^n [y_i - (\tilde{w}_j^T x_i + \tilde{w}_0)]^2$$

f_1, f_2, \dots
(no higher order variables)





M₁: overfit to
D_{Tr}



simplest [underfit]

DTR
100 H

DTe
120 H

{ perfectly →
tradeoff

20 Low

30 Low

Complex ✓

[overfit] →

0.1
v-low

50
not small

plots → practice

DTY

✓ (GOLD)

V-L
0-1

DTE

V-L
0-1

generalization

Malk

{ overfit - Underfit
Tradeoff

$$y_i = w_0 + \boxed{w_1} x_{i1}^1 + \boxed{w_2} x_{i1}^2 + \boxed{\sim w_3} x_{i1}^3 + \boxed{\sim w_4} x_{i1}^4$$

\dagger
(original)

Case 1: $w_1, w_2, w_3, w_4 \neq 0$

transformed feat uses
underfit

M3

$$w_3 = w_4 = 0$$

$$w_1, w_2 \neq 0$$

$$\begin{aligned} w_2 &= w_3 \\ &= w_4 = 0 \\ w_1 &\neq 0 \end{aligned}$$

Math:

weights on higher order poly $\not\rightarrow$

↳ more chance of overfitting

$$\{ R^2 \}$$

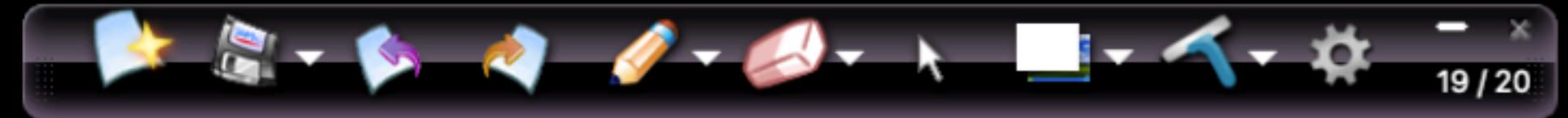
↑
Performance Metric

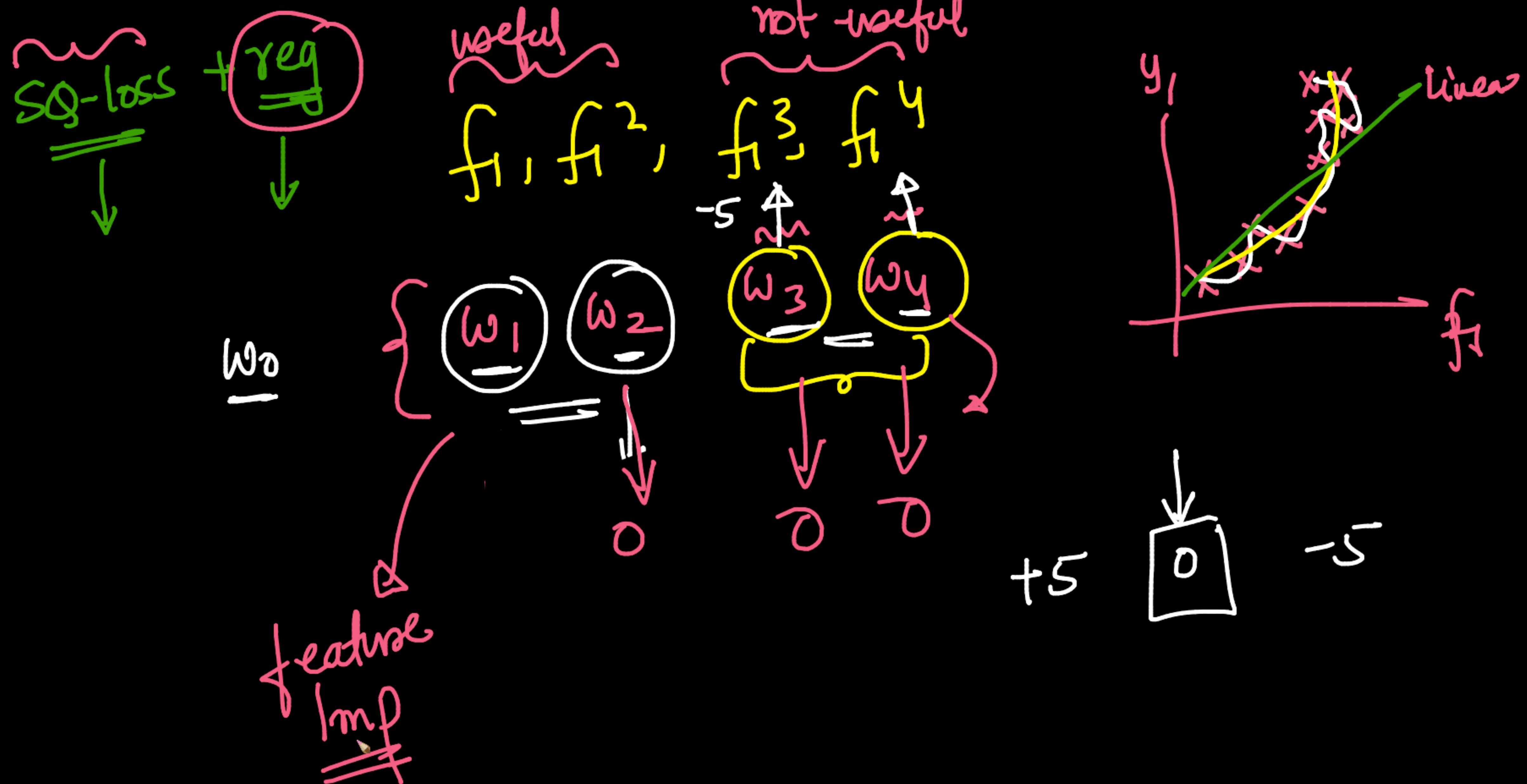
squared-loss

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Loss \rightarrow Minimize

perf-metric \rightarrow Maximize





Opt:

$$\min_{\underline{w}} \left[\sum_{i=1}^n \left(y_i - (\underline{w}^\top \underline{x}_i + w_0) \right)^2 + \lambda \sum_j w_j^2 \right]$$

square-loss

$$\min \left[\text{square-loss} + \lambda (w_j^s)^2 \right]$$

$\rightarrow f_1, f_1^2, f_1^3 \dots$

≥ 0

regularization

$$\min_{w_j} \left[C + \sum_{j=1}^d w_j^2 \right]$$

$w_j = 1 - \gamma_j$

$C + d \cdot 1$

$w_j^* = 0$

$w_j = 2$
 \downarrow
 $C + 4d$

$$\min_{w_j} \left[f(w_j) + \lambda \cdot \sum_{j=1}^d w_j^2 \right]$$

The diagram illustrates the optimization problem. A bracket on the left groups the function $f(w_j)$ and the regularization term. A bracket on the right groups the regularization term $\lambda \cdot \sum_{j=1}^d w_j^2$. A red arrow points from the variable w_j down to the term w_j^2 , indicating that the weight w_j is being squared. A red curly brace labeled "tve" is placed above the term $\lambda \cdot \sum_{j=1}^d w_j^2$, likely referring to a tuning parameter or a total variation penalty. An arrow points downwards from the right side of the equation towards the limit $w_j \rightarrow 0$.

$$\min_i \text{ Squared Loss}_i + \lambda \sum_j w_j^2$$

GD

$$\frac{\partial \mathcal{L}}{\partial w_j}$$

\mathcal{L}

$\min_{w_j} \sum_i \text{SQ. loss}_i + \lambda \sum_j w_j^2$

Q

$\lambda \sum_j |w_j| \rightarrow \text{reg}$

$L_1\text{-reg}$

$w^T w = \|w\|_2 = L_2\text{ norm}$

* not diff @ $w_j = 0$

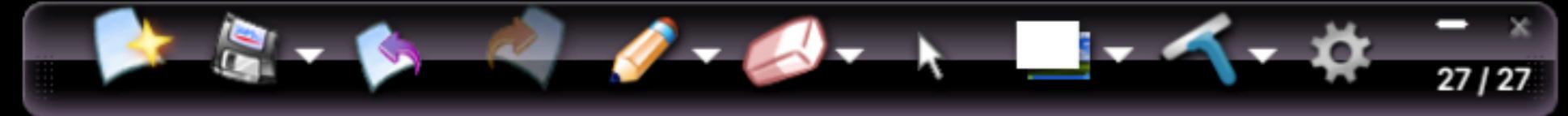
hack (later)

w_j

+ve, -ve

+ve, -ve

post class → revise (think & revise)



$$\begin{matrix} f_1 \\ w_1 \end{matrix} \quad \begin{matrix} f_1^2 \\ w_2 \end{matrix} \quad \begin{matrix} 2f_1^2 \\ w_3 \end{matrix}$$

$$5 \quad \boxed{0} \cdot \boxed{0}$$

$$5 \quad \boxed{\overline{Z}} \quad \boxed{\overline{Z}}$$



$$y_i = \underline{w^T x_i} + w_0$$

parameters

min $\sum_{i=1}^n (y_i - \hat{y}_i)^2$ SQ. loss

\hookrightarrow 10000

λ reg

$\boxed{\lambda}$ hyperparam

\boxed{LB}

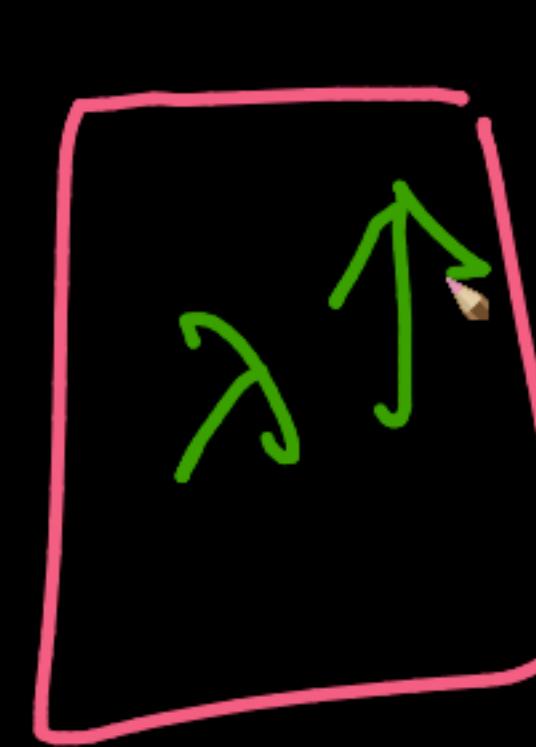
\rightarrow simpler

{ Case 1 :
Case 2 :
Case 3 :
Case 4 :

$\left\{ \begin{array}{l} \lambda = 0 \Rightarrow \text{no reg} \rightarrow \text{overfit} \\ \lambda = 10 \\ \lambda = 100 \\ \lambda = 10000 \\ \lambda = \text{BigNum} \Rightarrow \infty \end{array} \right.$

Underfit





overfit → underfit
↑
good-model

$$\min_{\mathbf{w}_j} \sum_{i=1}^n [y_i - (\mathbf{w}^\top \mathbf{x}_i + w_0)]^2 + \lambda \sum_j w_j^2$$

GD on each w_j

$\lambda (\mathbf{w}^\top \mathbf{w} - 1)$

PCA

Constraint opt

objective

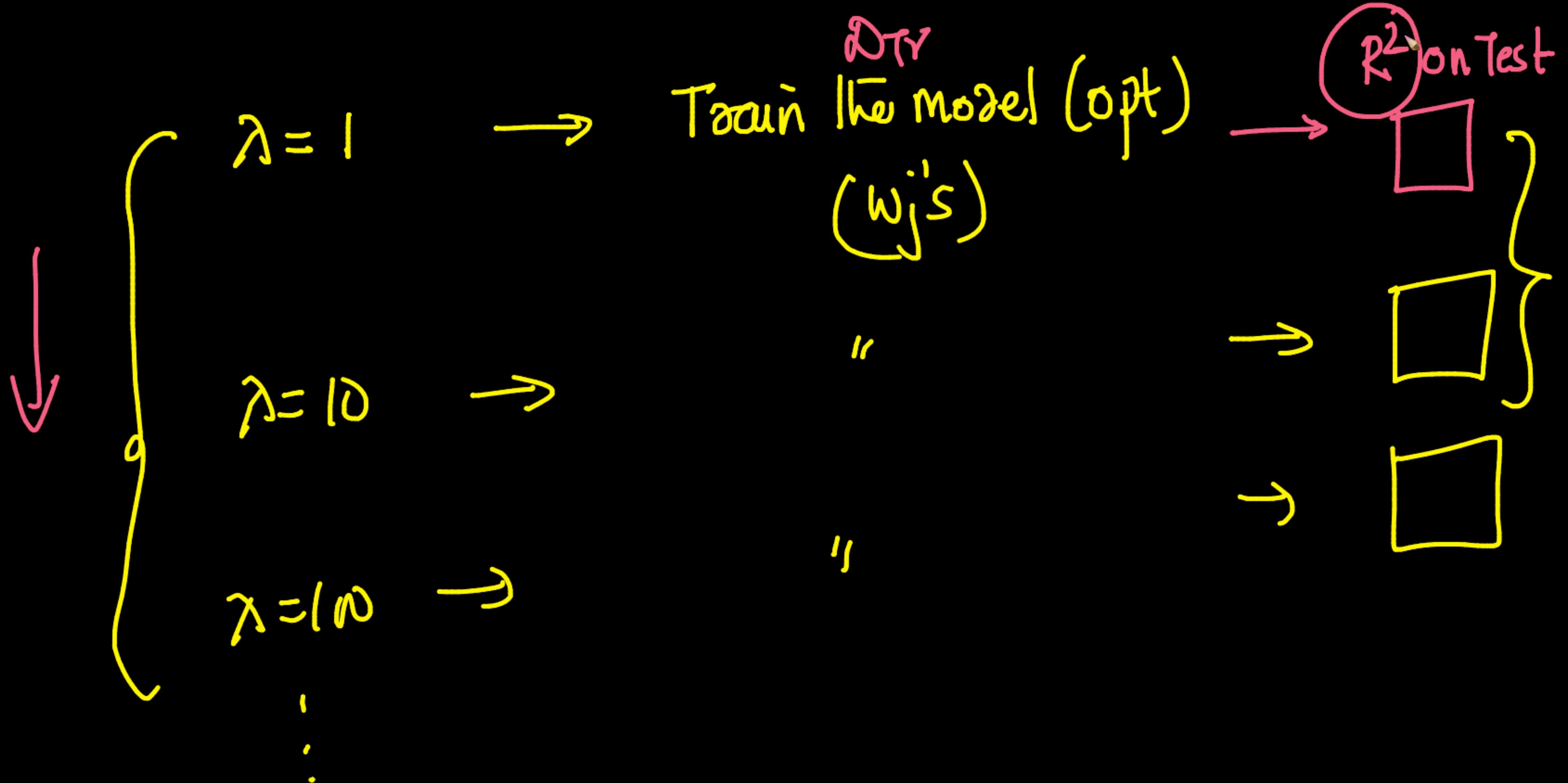
Constrained optimizn:

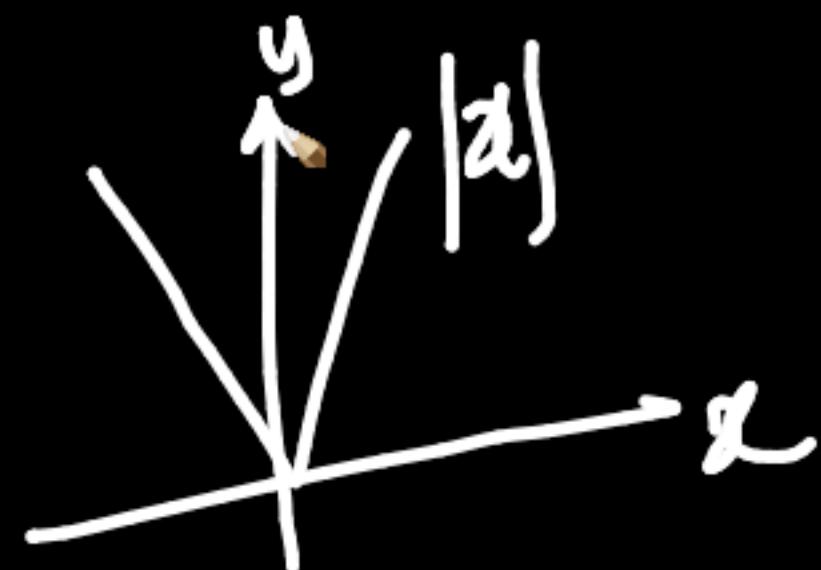
$$\frac{\partial \mathcal{L}}{\partial \lambda} = 0$$

GD on λ
to find $\boxed{x^*}$

→ Yes

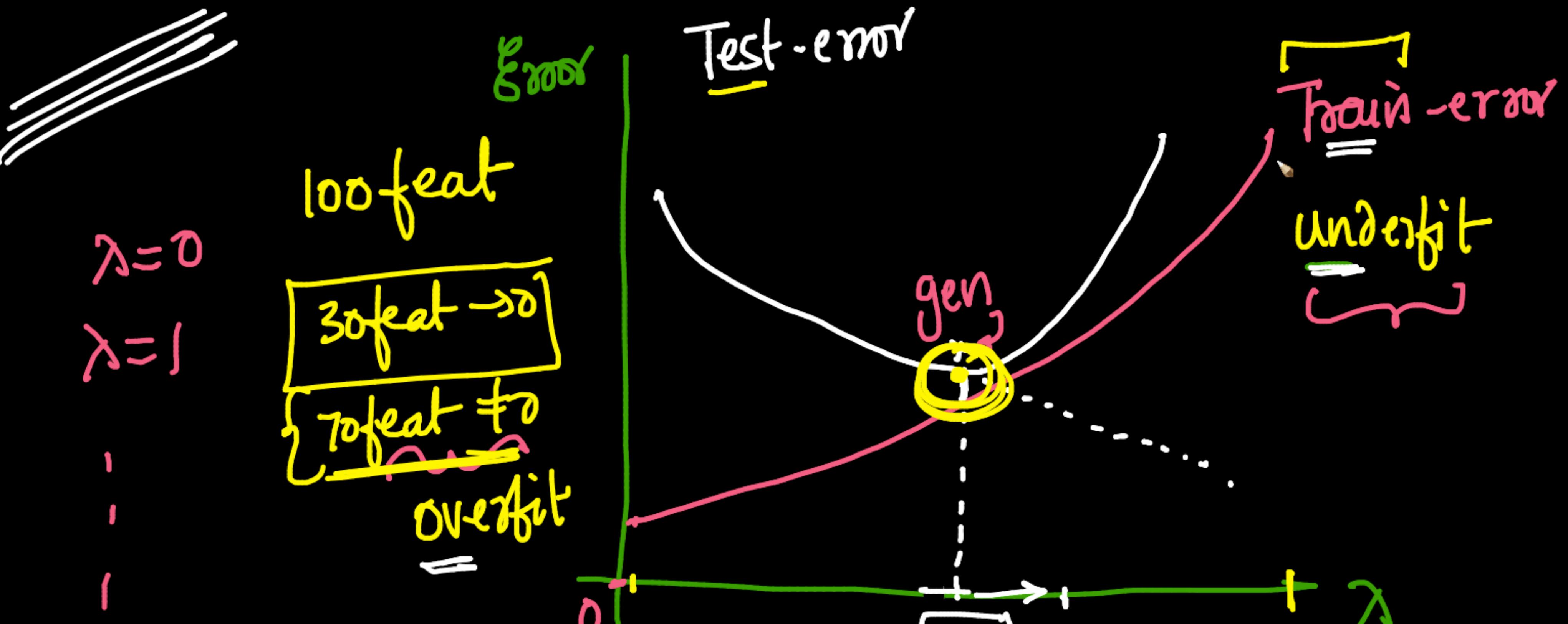
ML-optimizn: No
(manual) → $\boxed{\lambda}$ control [tradeoff]
overfitting — underfit





$$\mathcal{L} \leftarrow \min_{w_j} \sum_{i=1}^n \left[y_i - (\underline{w}^T x_i + w_0) \right]^2 + \lambda \sum_j w_j^2$$

$\left\{ j \neq 0; \frac{\partial \mathcal{L}}{\partial w_j} = \sum_{i=1}^n 2(y_i - (\underline{w}^T x_i + w_0)) (-x_{ij}) + \lambda \cdot 2 \cdot w_j \right\}$



$$\min_{w_j} \sum_i (y_i - \hat{y}_i)^2 + \lambda \sum_j w_j^2$$

100 features \rightarrow 70feat fo ; 30 \Rightarrow some slightly useful features

$$\min \text{loss} + \lambda_B \min \text{weights}$$

$$\lambda_B = 10$$

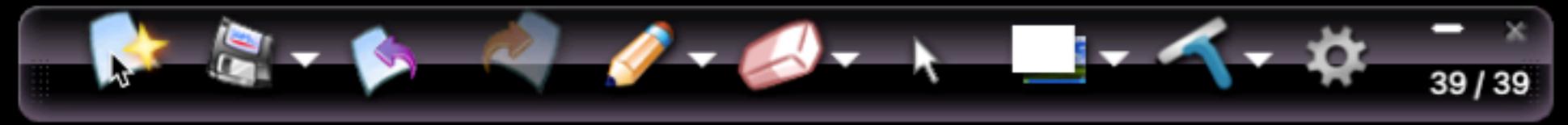


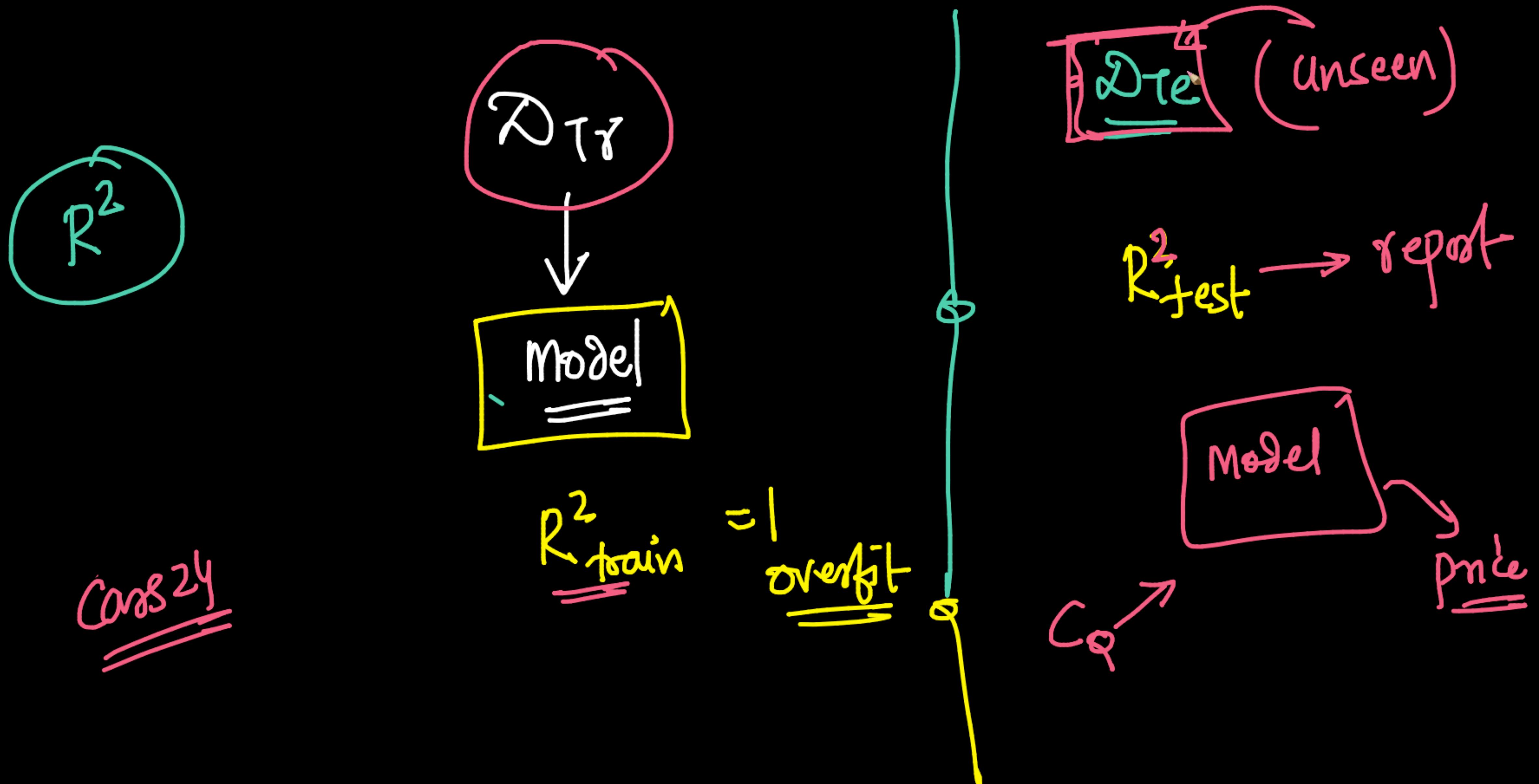
Test-data
unseen



$\lambda = 1 \rightarrow$ Train the model
 $= \bar{w_j} =$

$\lambda = 2 \rightarrow$ " $= \bar{w_j} =$





Test - error

100

102

non-zero
Weights

mm-weights

$\lambda = 10$

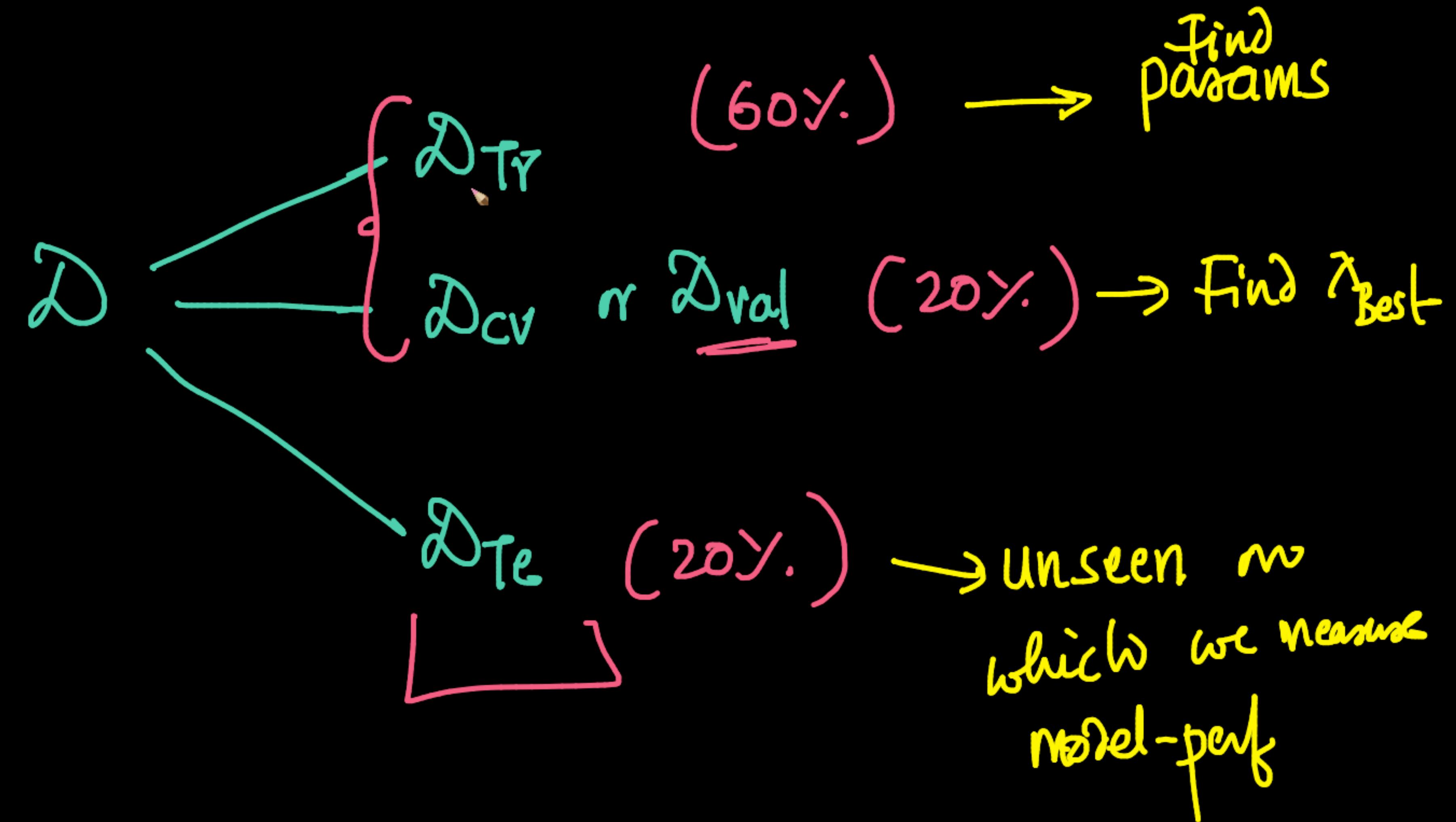
$\lambda = 9$

Model-Simplicity vs Model perf

{
explanability

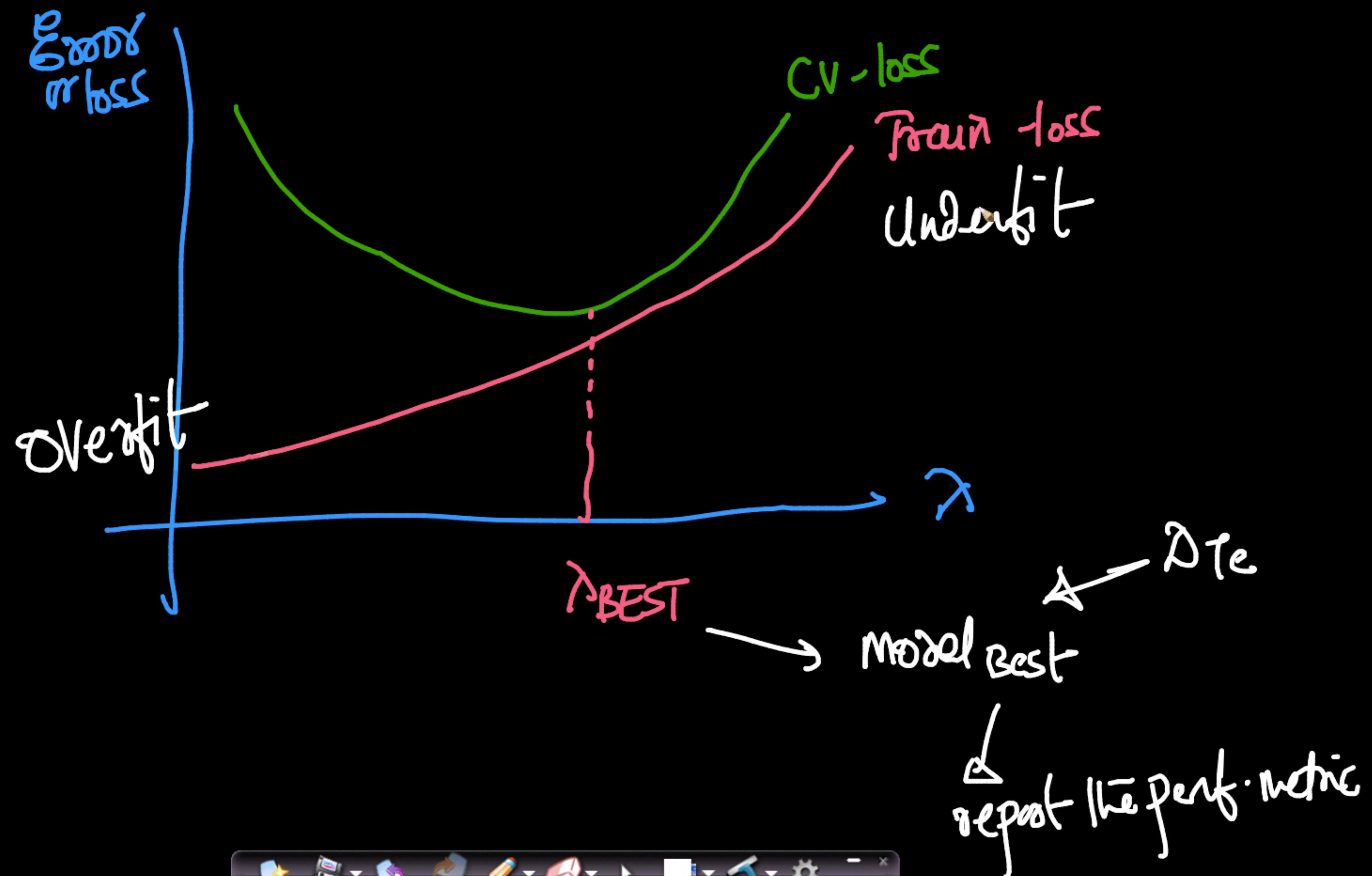
✓ Using test-data to find $\hat{\gamma}_{\text{BEST}}$ is not correct

Soln:



did
not
use
 D_{Te}

- For $\lambda = 1, 10, 20, 50, 100, \dots$
 - Train a model m on \mathcal{D}_{TR} with λ fixed
 - $(w_j^s) \leftarrow \overline{J}$
 - Measure model-perf (R^2) on \mathcal{D}_{CV}
- Pick the λ_{Best} s.t. it has lowest CV-error
 - $\rightarrow m_{ModelBest}$
- Measure the perf of $m_{ModelBest}$ on D_{Te} & report it



500 p|s

300
100
10

features \rightarrow 7

$d \ll n$

K-fold CV

Overfitting - Underfitting tradeoff

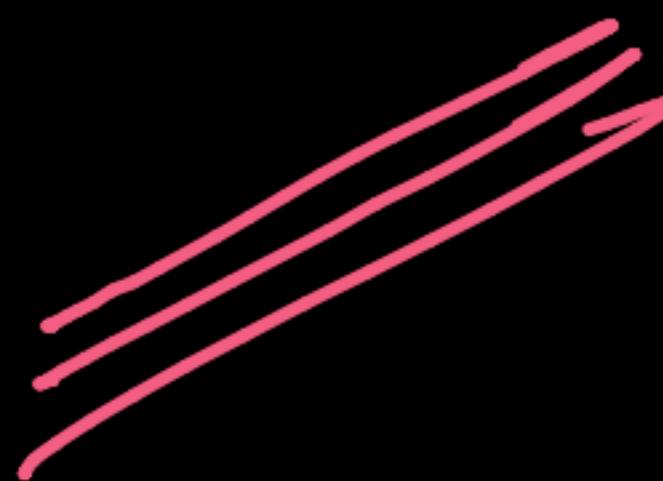
high Variance - Bias tradeoff

$$y_i - \hat{y}_i = e_i$$

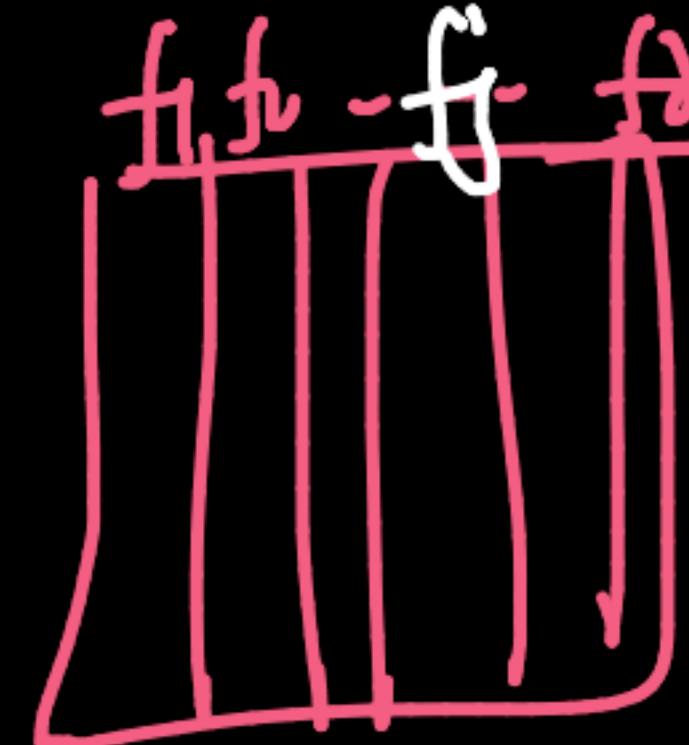
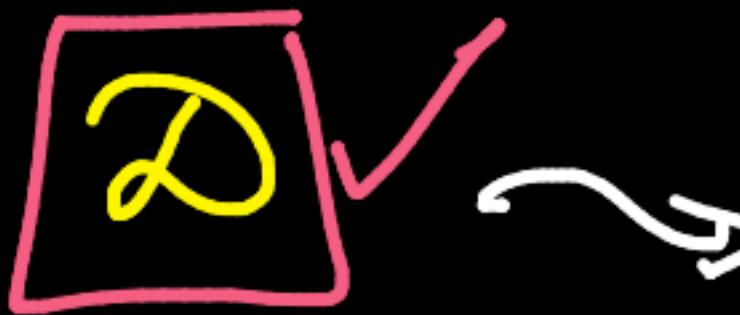
labeled $y_i - \hat{y}_i$ → e_i (Bias + Variance) \rightarrow Expectations $E(x)$

labeled $y_i - \hat{y}_i$ → e_i (Bias + Variance) \rightarrow Expectations $E(x)$

labeled $y_i - \hat{y}_i$ → e_i (Bias + Variance) \rightarrow Expectations $E(x)$



Summarize



①



\bar{D}_{TY} , \bar{D}_{CV} , \bar{D}_{Pe}

✓

②

col. standardize

\bar{D}_{TY}

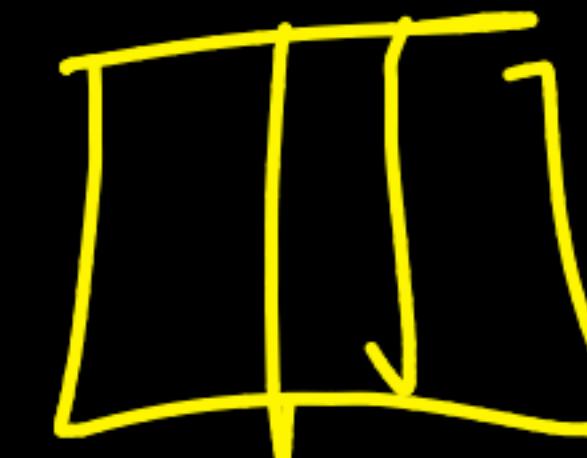
$\forall j: m_j; s_j$ (Shoe)

③

Using m_j & s_j $\rightarrow \bar{D}_{CV}$ & \bar{D}_{Pe}

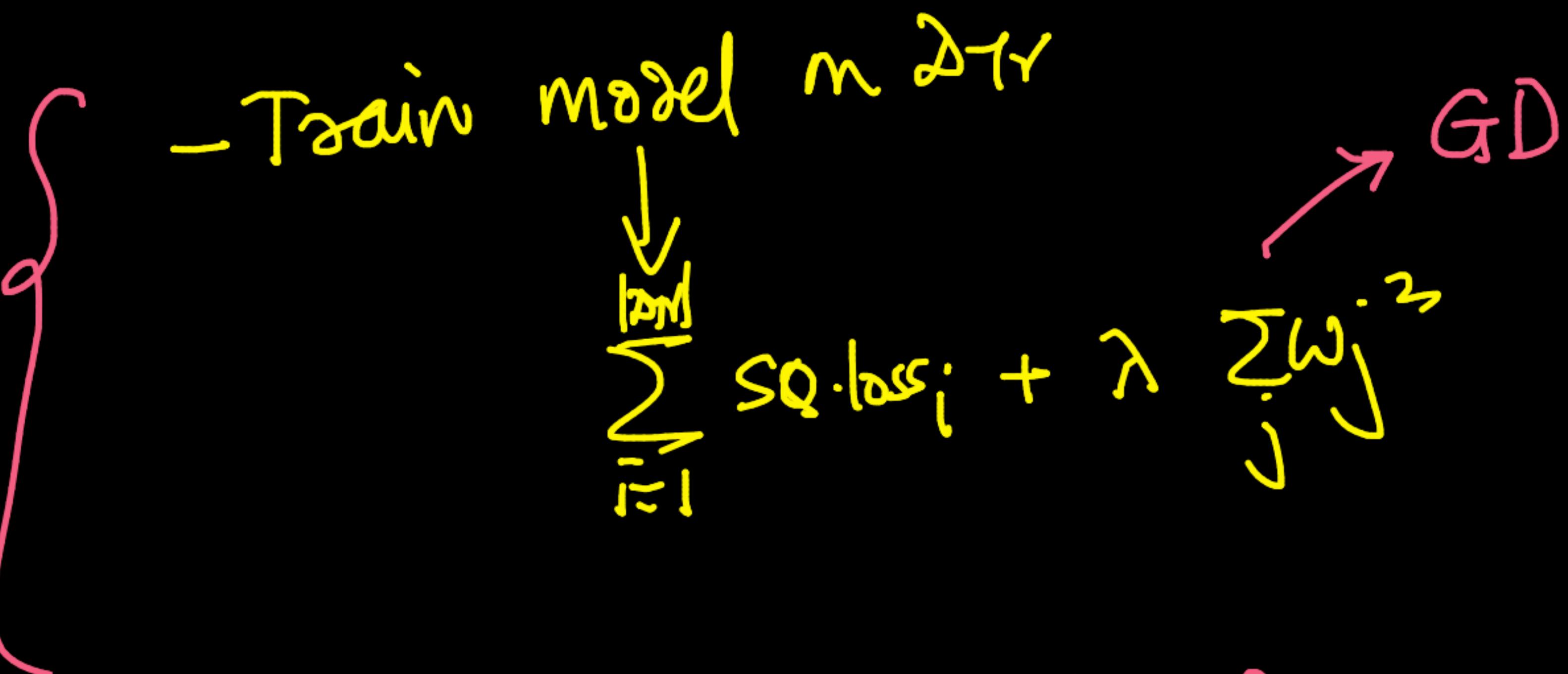
Transform

DK



④

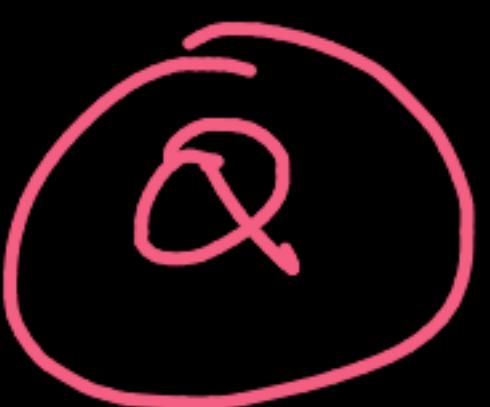
For $\lambda = 1, 10, 20, \dots$



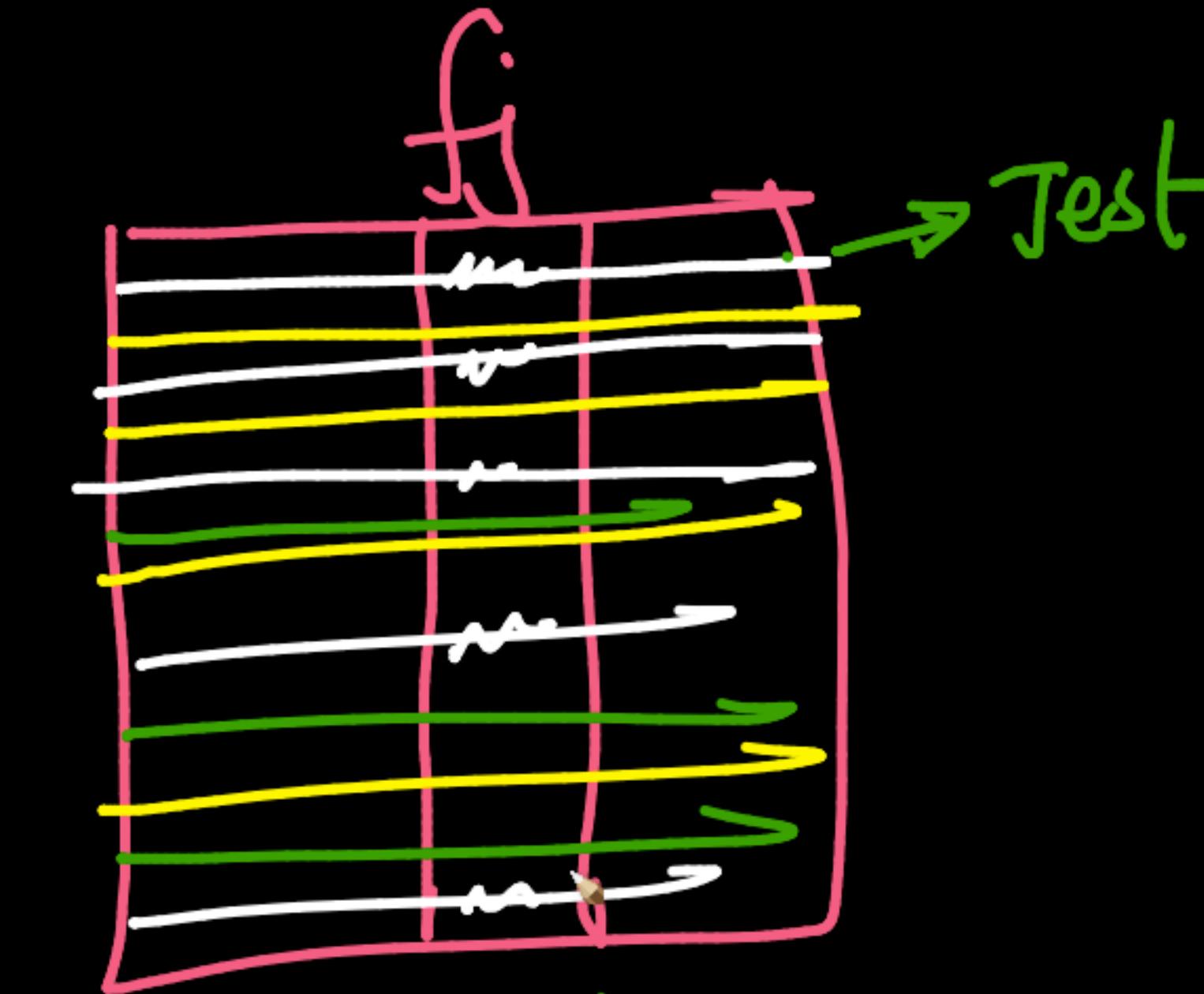
⑤

pick λ_{BEST} using error or perf. on \mathcal{D}_{CV}

⑥ Report Model-perf using ModelBest on TestJdl-



2



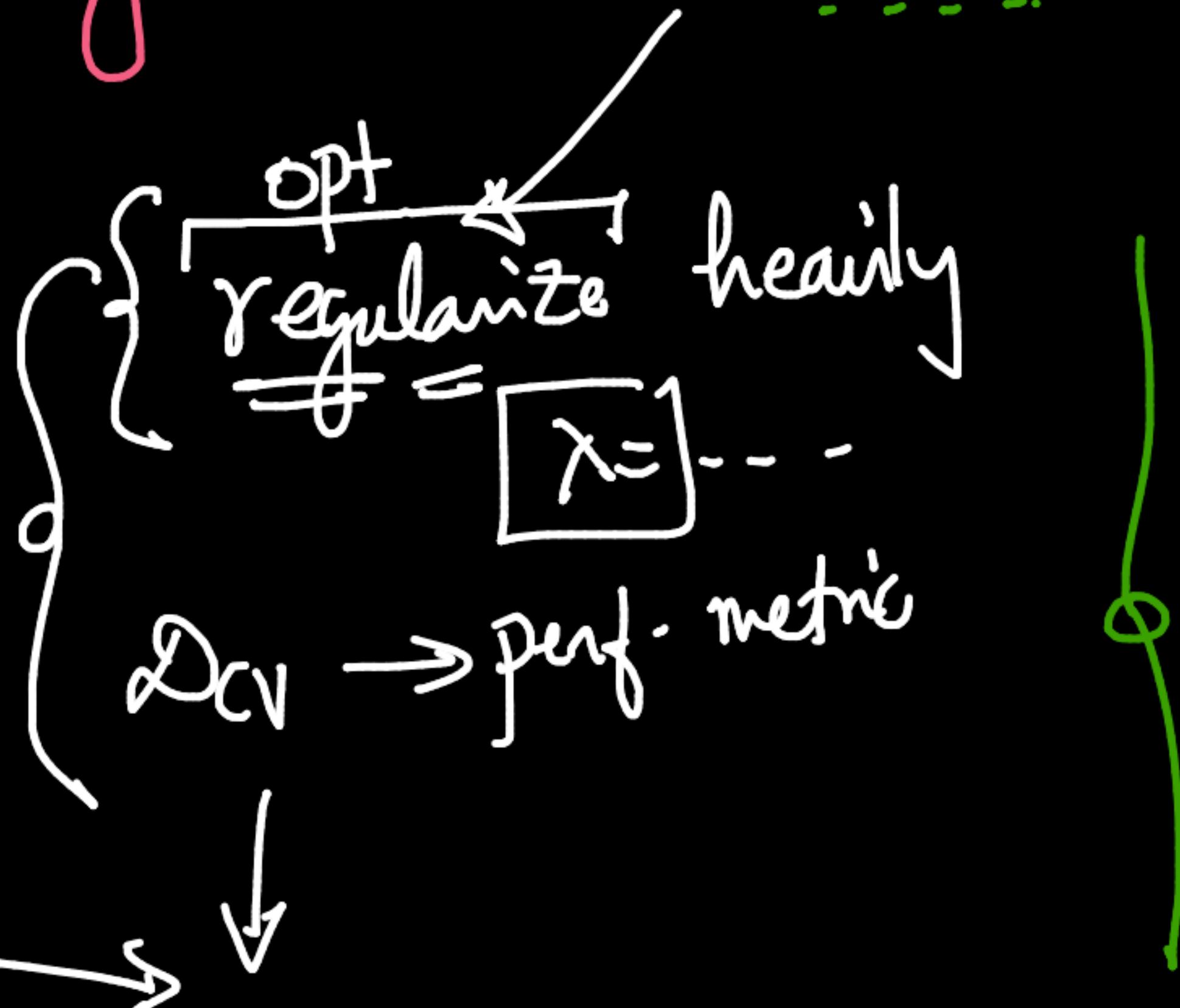
Leakage
~~if~~

↳ Some info about
Die leakage into
your邦定 process

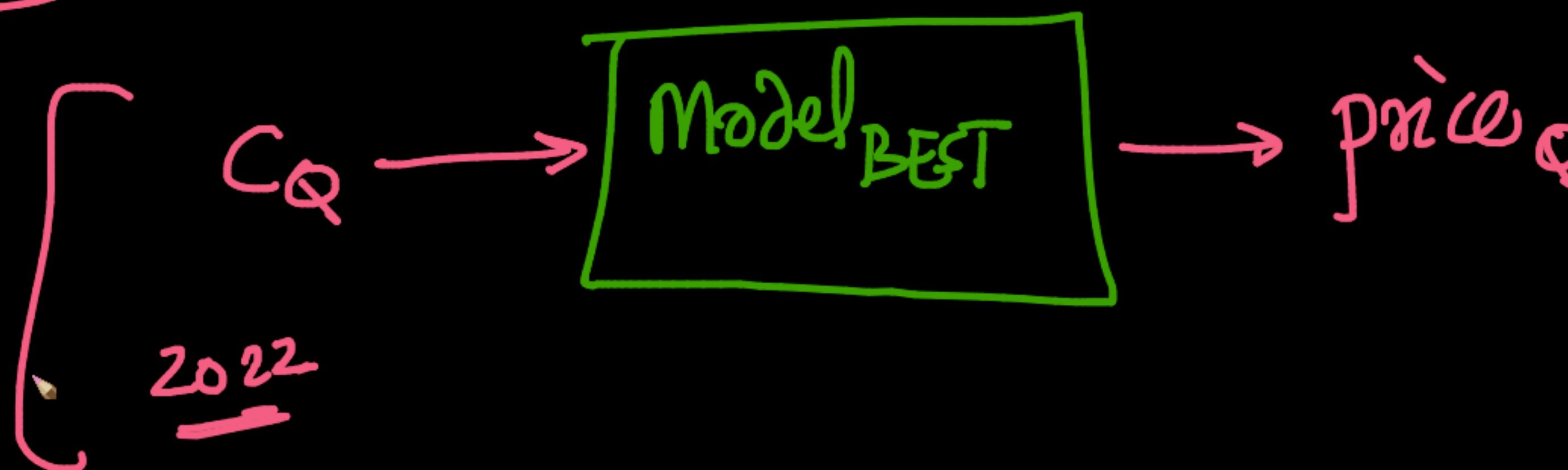
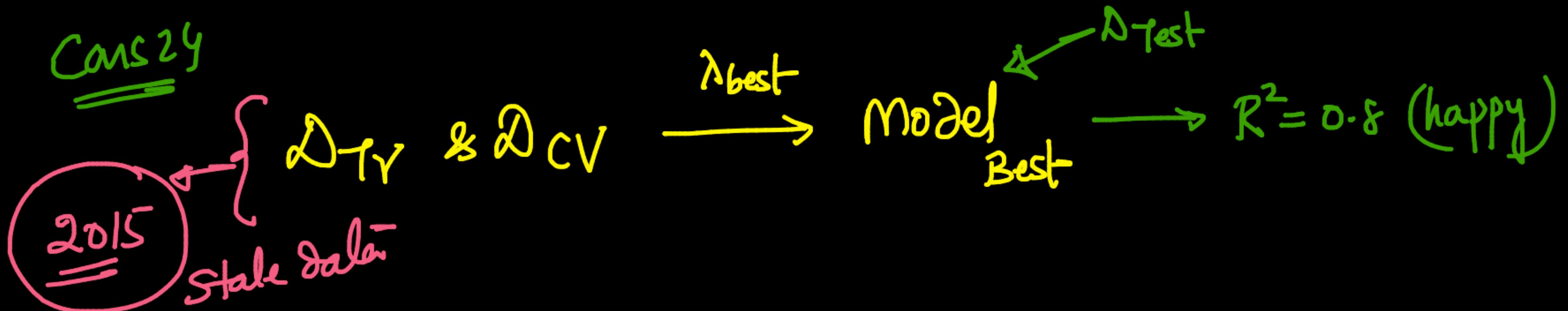
↳ ~~m_j & S_j~~ ↳ DTR ↳ Model
params



L1-reg → lots of features → More chances of useless or correlated features



stats: - VIF to get rid of correlated
- EDA & plotting



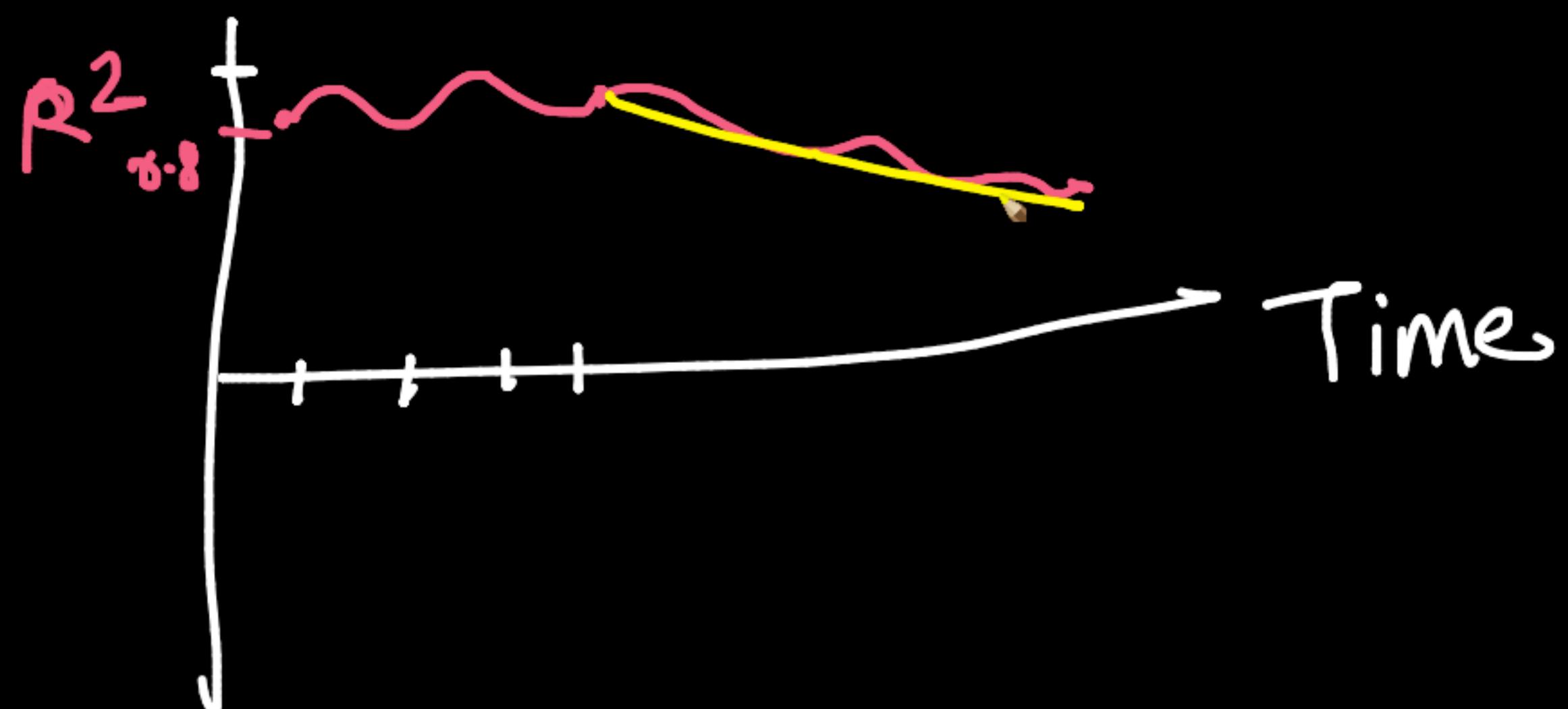
✓ { - retrain my model every month on recent data

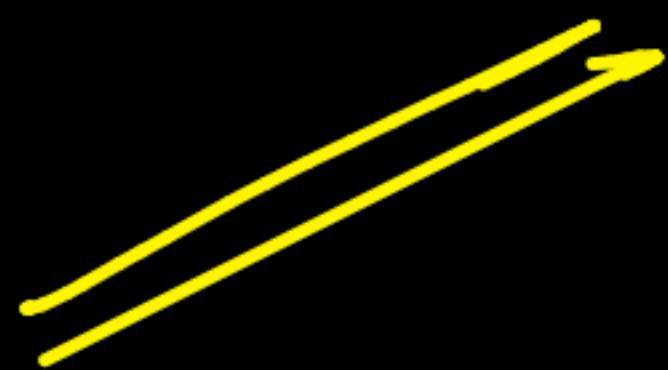
-

Model is deteriorating over time

$$D_{\text{Test}} \rightarrow R^2 = 0.8$$

$$\begin{aligned} & \text{pred. price} - \text{act price} \\ & \hat{y}_i - y_i \end{aligned}$$





d-features

Model Best

-



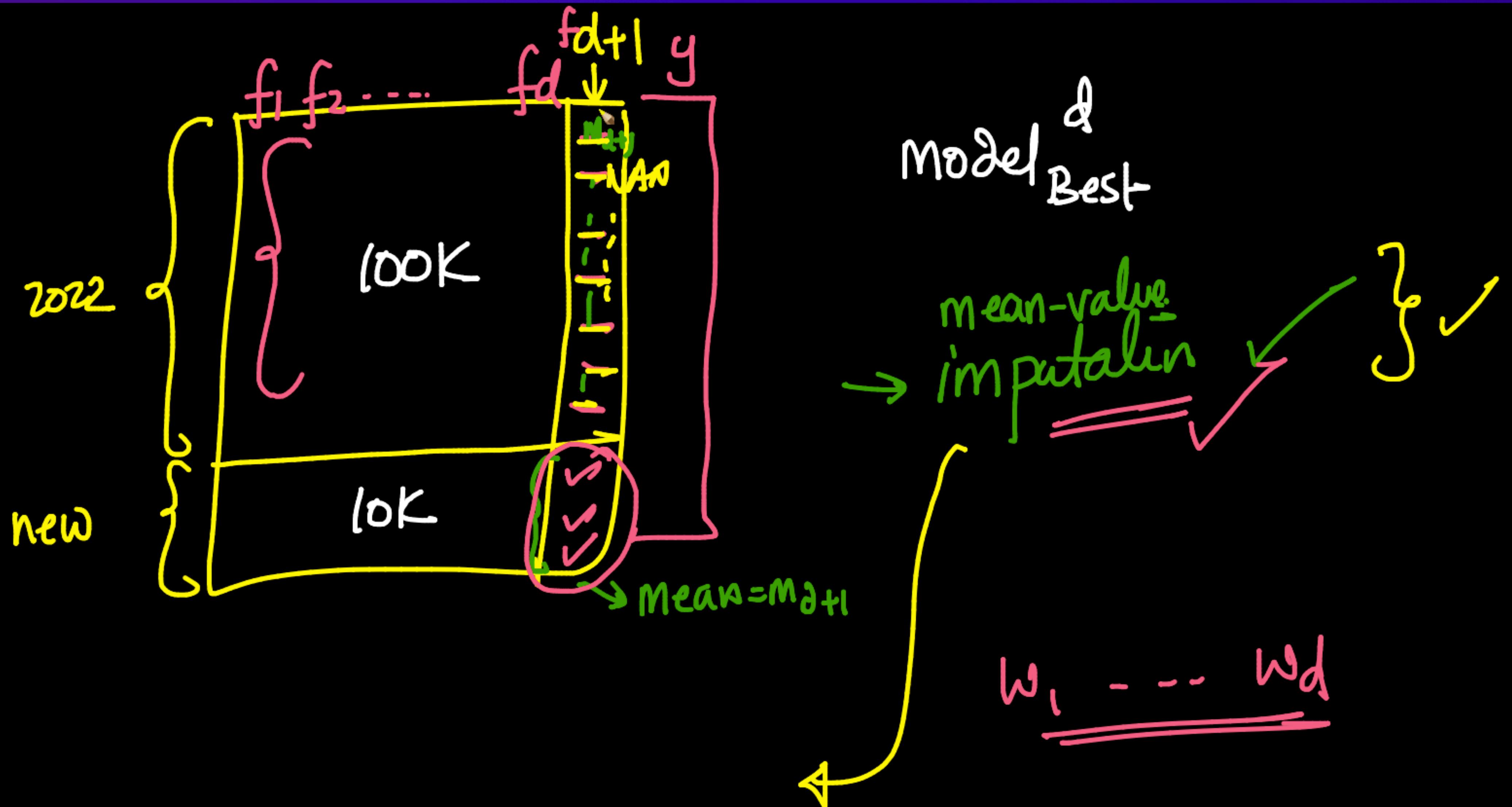
2022

d+1st features

- VIF of 2+1st feat with exist^{ing} d-feat
- retain with all d+1st features



Model perf
improve



opt-hack

look

$\checkmark \frac{r_0}{w_j} + \frac{q}{w_j} x_{i,1} + \dots + \underbrace{\frac{w_d}{w_j} x_{i,d}}_{\text{up}} + w_{d+1} x_{i,d+1}$

var

$$\min_{w_{d+1}} \sum_i \text{loss}_i + D \sum_j w_j^2$$

$$\frac{\partial L}{\partial w_{d+1}}$$

GD

w_{d+1}

{revise & think deeply}



_ Code (trivial)

- logistic Reg.

Test error

100

 $\lambda = 10$ weights

8 non-zero params = 2 zero params

 $\lambda = 20$

fewer 3 non-zero params = 7 zero params

~~120~~ ✓

Fewer non-zero feat \Rightarrow few comp

- Computational
- Battery
- Memory
- Interpretability (Simpler)
/Explanability

