

FoDA - L13

Cross - Validation

How well is regression working?

Input $(x, y) \quad x, y \in \mathbb{R}^n$

convert $x \rightarrow \tilde{X}_P = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^P \\ 1 & x_2 & x_2^2 & \dots & x_2^P \\ \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix}$

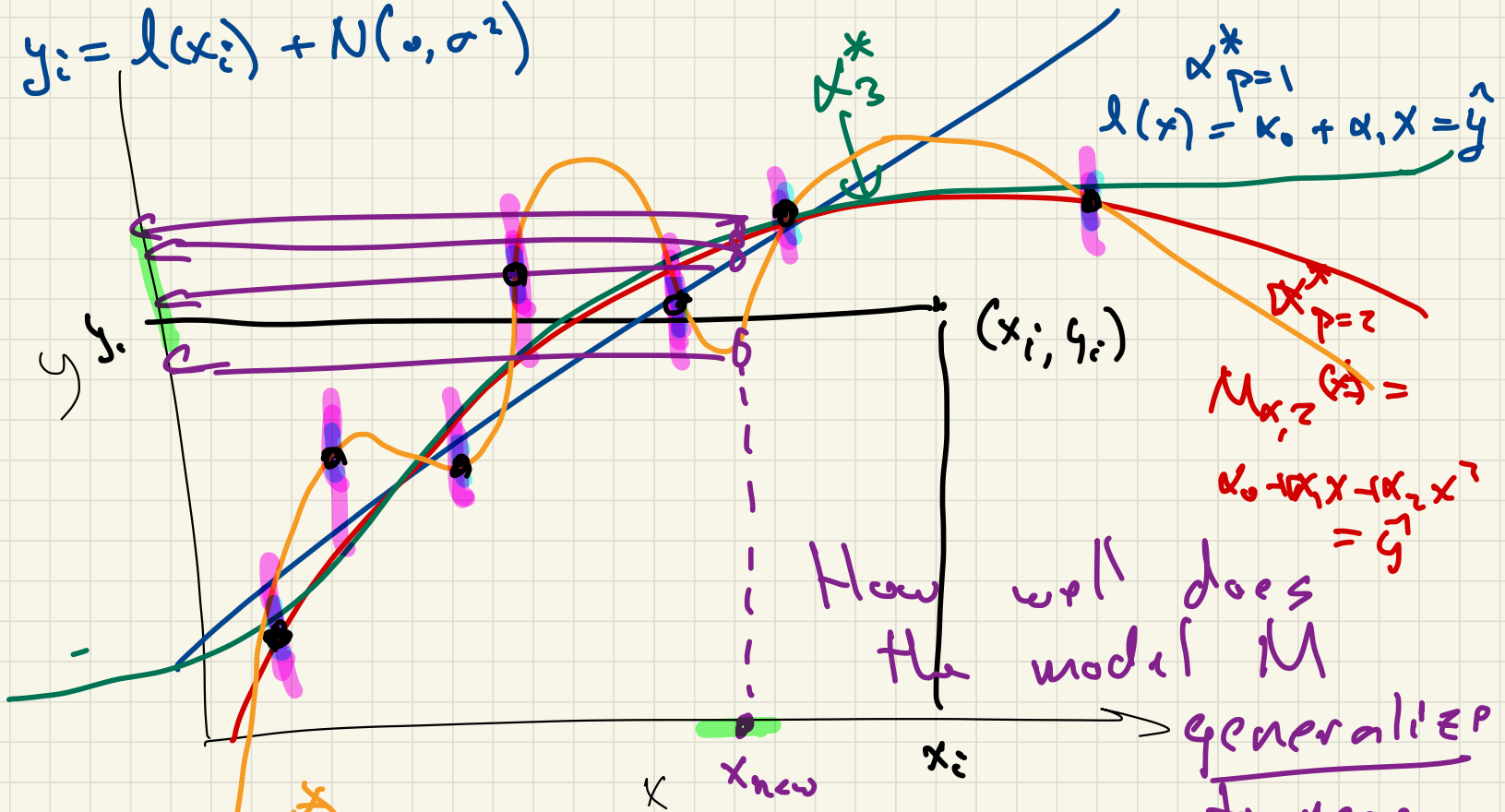
$$\alpha_P^* = (\tilde{X}_P^T \tilde{X}_P)^{-1} \tilde{X}_P^T y$$

$$\alpha_P^* = (\alpha_0, \dots, \alpha_P) \in \mathbb{R}^{P+1}$$

$$M_{\alpha_{1P}}(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \dots + \alpha_P x^P$$

Which choice of P ?

$$y_i = l(x_i) + N(0, \sigma^2)$$



$$l(x) = \alpha_0 + \alpha_1 x = \hat{y}$$

$\alpha_1^* = 1$

$$\alpha_1^* = 2$$

$$M_{\alpha, 2}(\hat{y}) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 = \hat{y}$$

$$\alpha_p^* = 7$$

$$SSE((X, y), M_{\alpha^*, 7}) = 0$$

How well does the model M generalize to new data?

What makes a model good?

stable to noise (x-coord)

- fit data well (w/ noise tolerance in y)

simple as possible

↳ How well does it generalize to new data?

Cross-Validation

Assume

$$(X, y) \stackrel{\text{iid}}{\sim} \phi$$

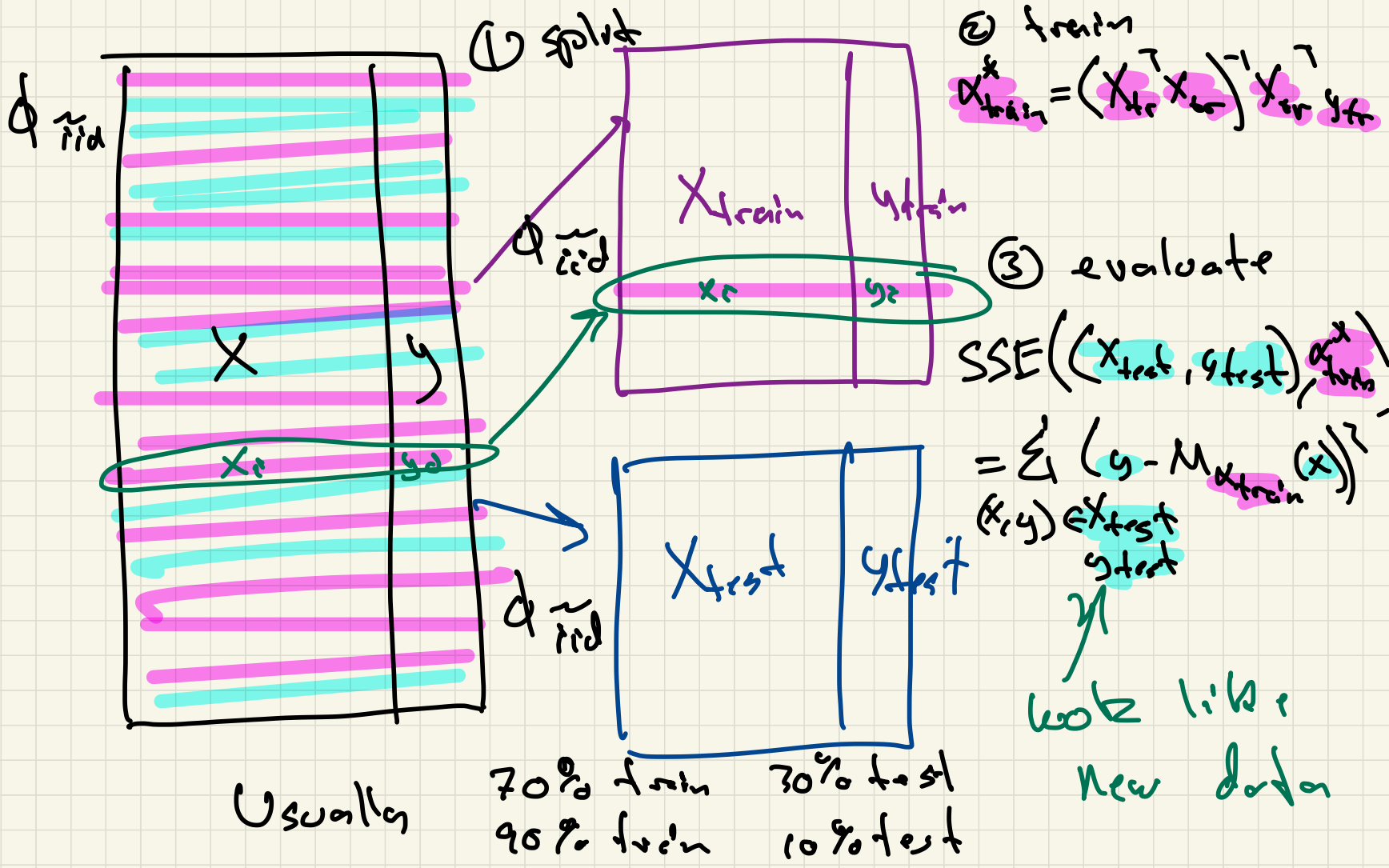
← only have
of input

New data $(x', y') \stackrel{\text{iid}}{\sim} \phi$

Assume n (# data points) is large

Randomly split (X, y) into training set
 $(X_{\text{train}}, y_{\text{train}})$

and a test set
 $(X_{\text{test}}, y_{\text{test}})$



What is cross-validation good for?

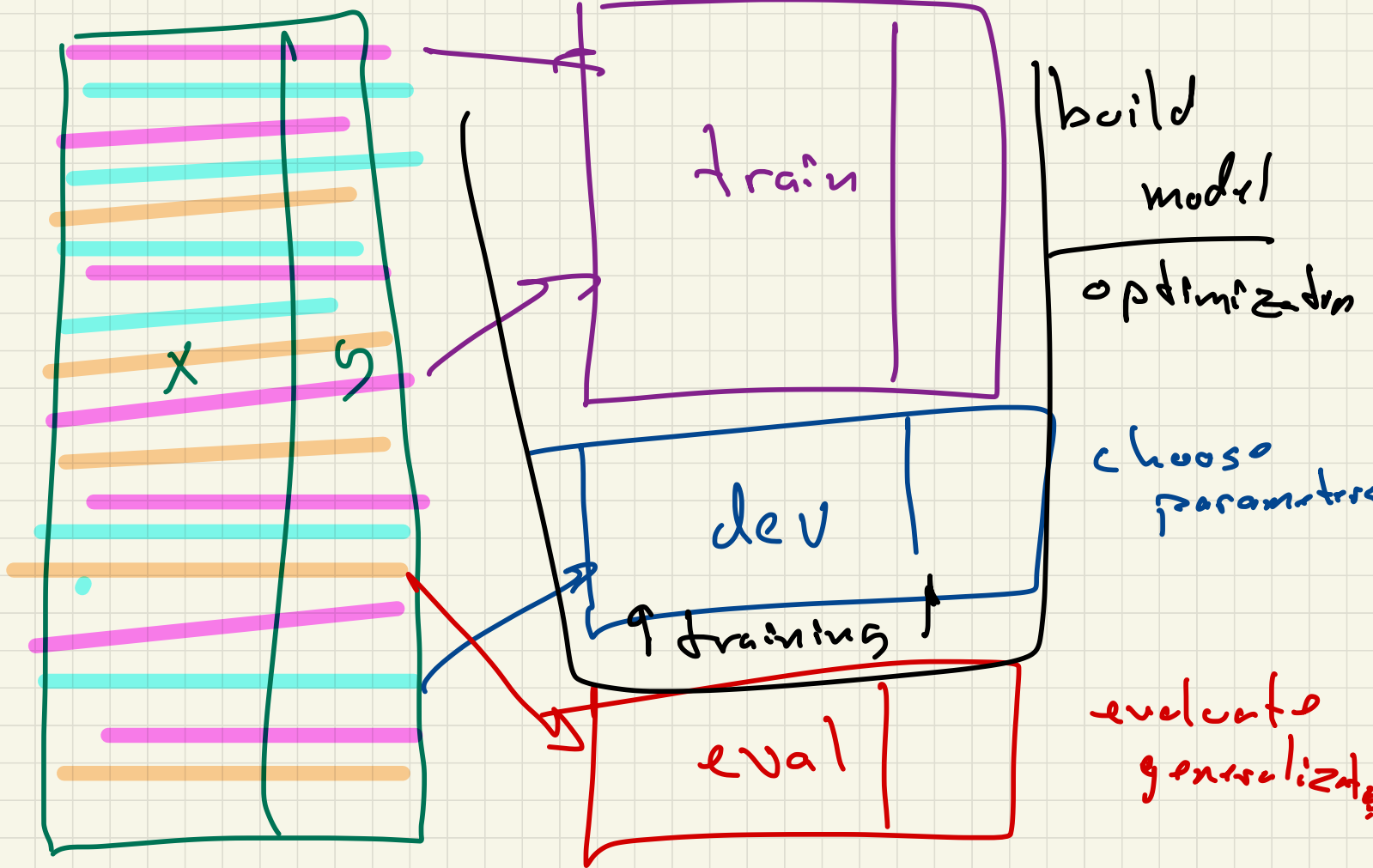
- ① Predict how well does M_α do on new data.

$$\text{RMSE} = \sqrt{\frac{1}{|X_{tr}|} \sum_{x \in X_{tr}} (y - M_\alpha(x))^2}$$

root mean sq. err

avr. sq. error
MSE

- ② Choose the best model (parameter)
- $\text{MSE}(X_{tr}, y_{tr}, M_{\alpha, 2})$ vs. $\text{MSE}(X_{tr}, y_{tr}, M_{\alpha, 7})$ e.g. **A**



$$p^* = \underset{\uparrow}{\text{arg min}} \quad \text{SSE}((x_{te}, y_{te}), M_{\alpha_p})$$

$$= \underset{\uparrow}{\text{arg min}} \quad \begin{aligned} &\text{RMSE}((x_{te}, y_{te}), M_{\alpha_p}) \\ &= \sqrt{\frac{1}{n} \text{SSE}} \end{aligned}$$

What to do w/ not
enough data to split?

Leave-one-out C-V n data points

Create n test / train splits

$X_{tr, i}$ $i \in [1, \dots, n]$ ($X_{tr, 1}, X_{tr, 2}, \dots, X_{tr, n}$)

$X_{tr, i} = \{x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}$

$X_{te, i} = \{x_i\}$

$$\text{Error} = \sum_{i=1}^n (y_i - M_{x_{\text{test},i}}(x_i))^2$$

Single
test
point
for
model

build model
 $M_{x_{\text{test},i}}$ for each
split.