



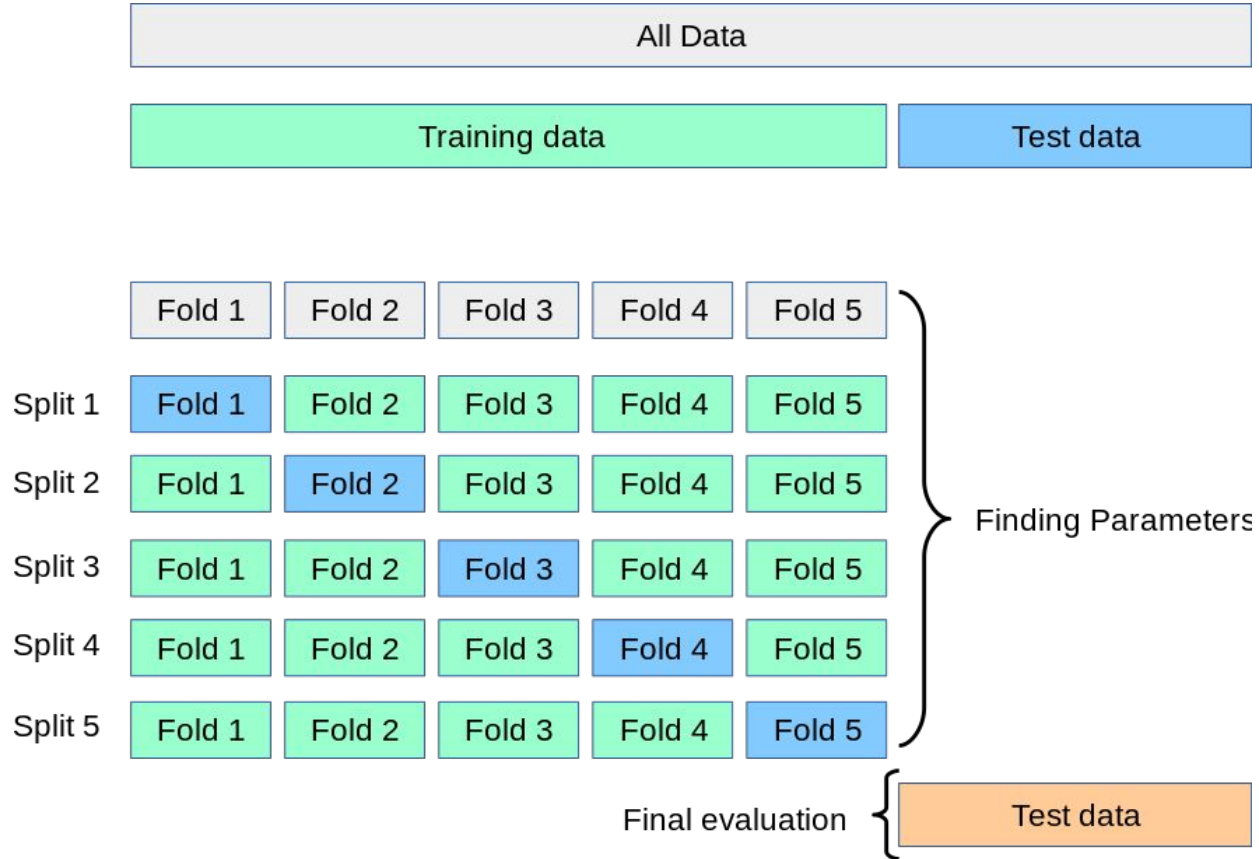
# Lecture 8: Bias-Variance tradeoff; More Ensembling

# Outline

1. Validation Strategies
2. Blending
3. Stacking
4. Bias-Variance Tradeoff

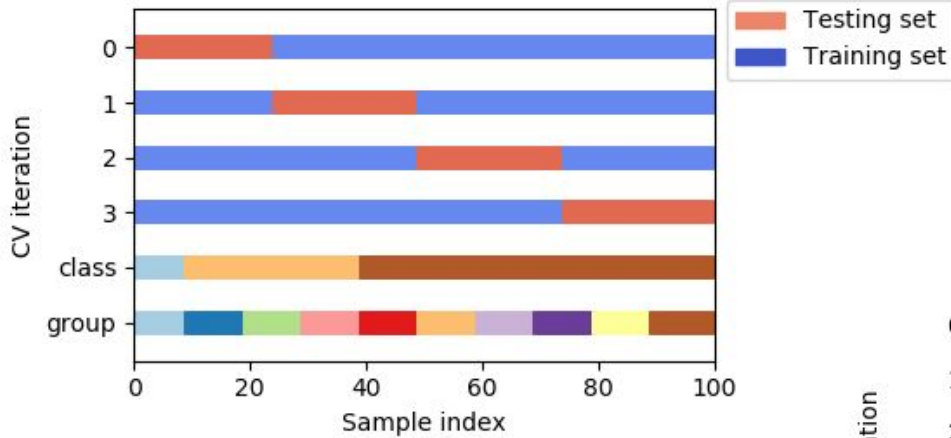
# Validation Strategies

# Validation strategies

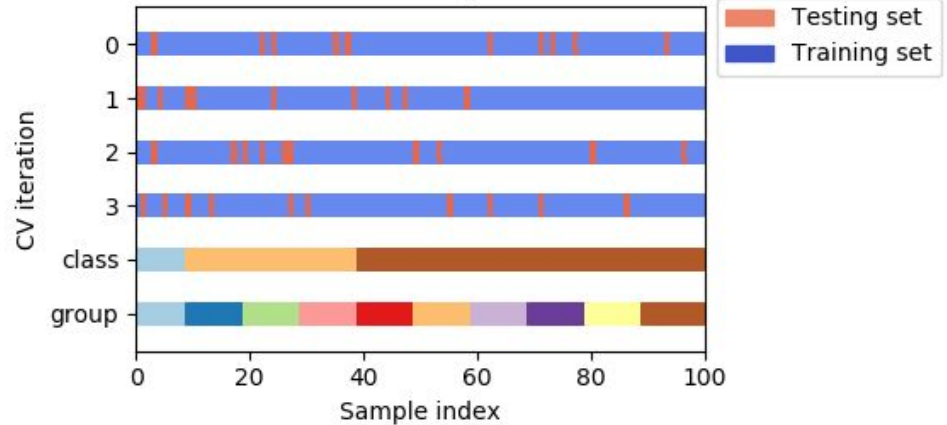


# Validation strategies

KFold

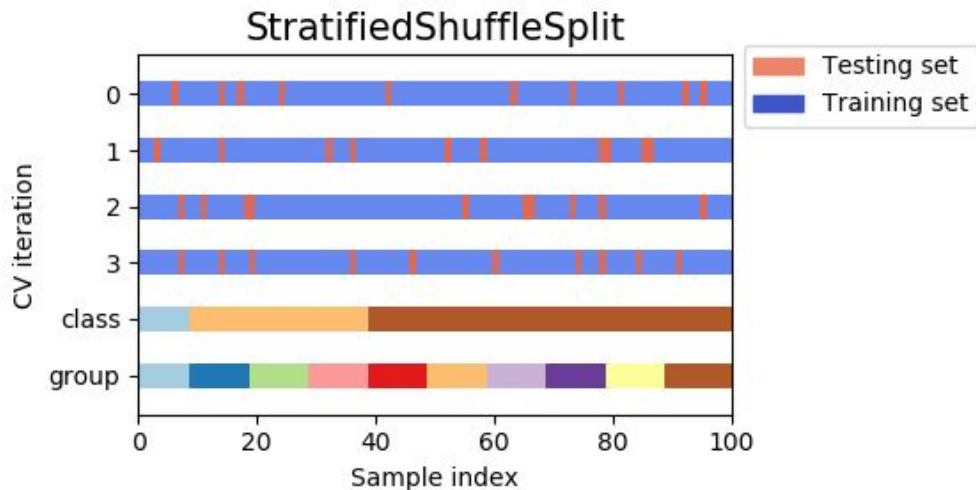
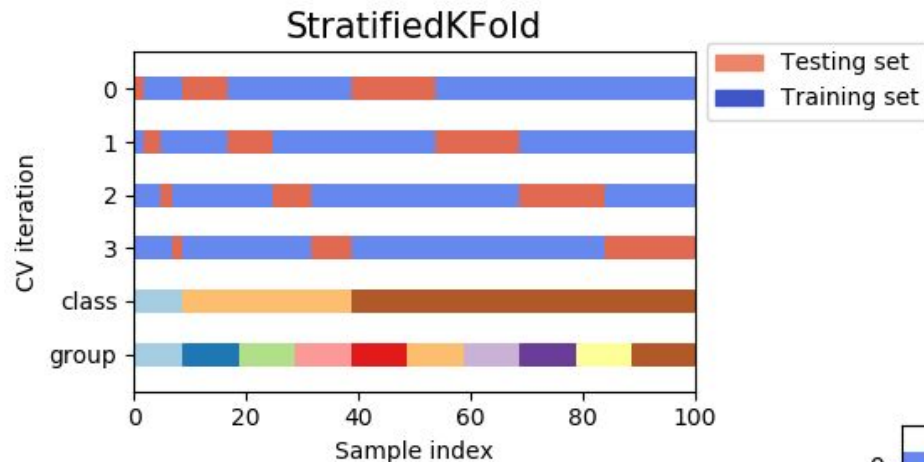


ShuffleSplit

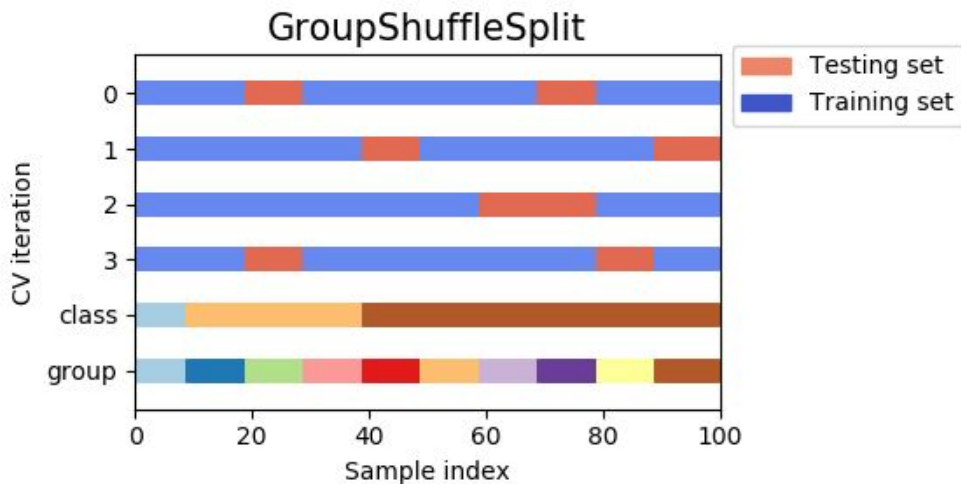
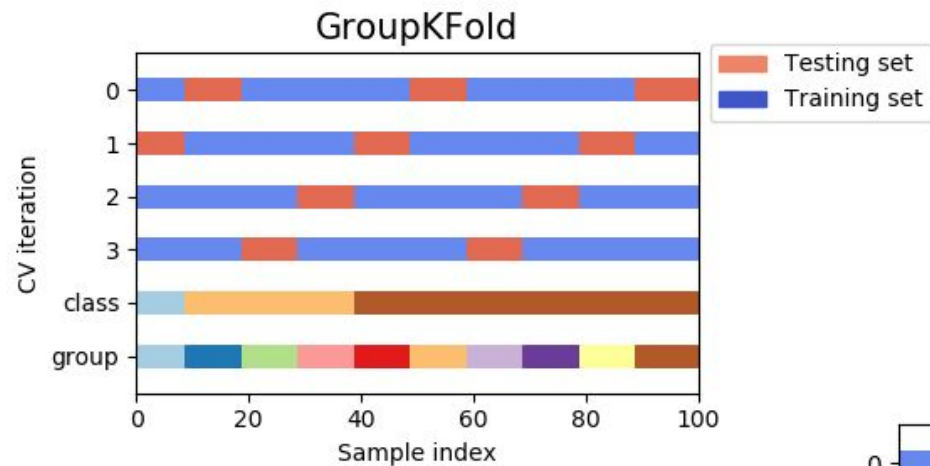


Special case: Leave One Out (LOO) - good for tiny datasets

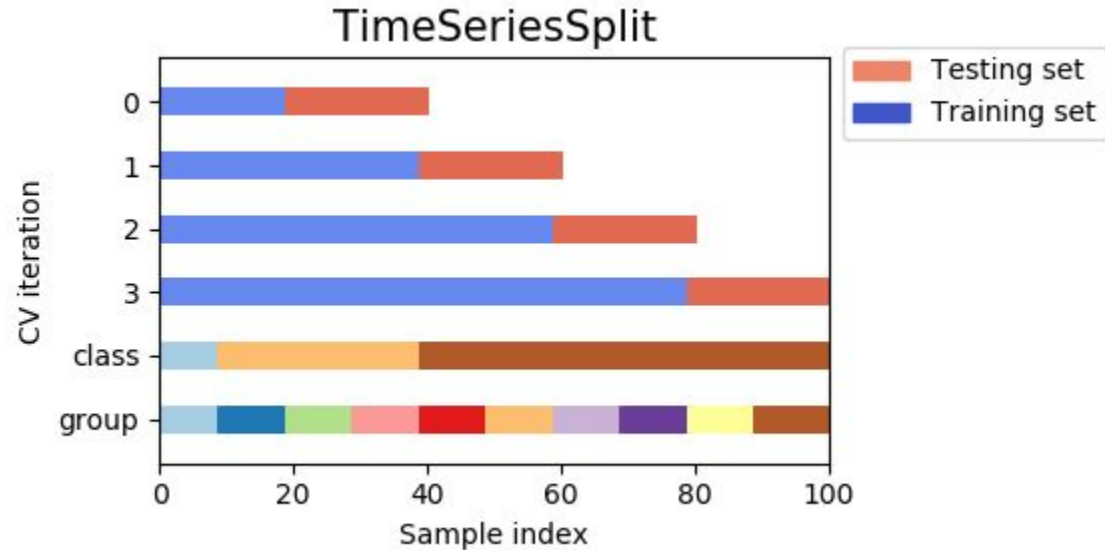
# Validation strategies



# Validation strategies



# Special case: time series



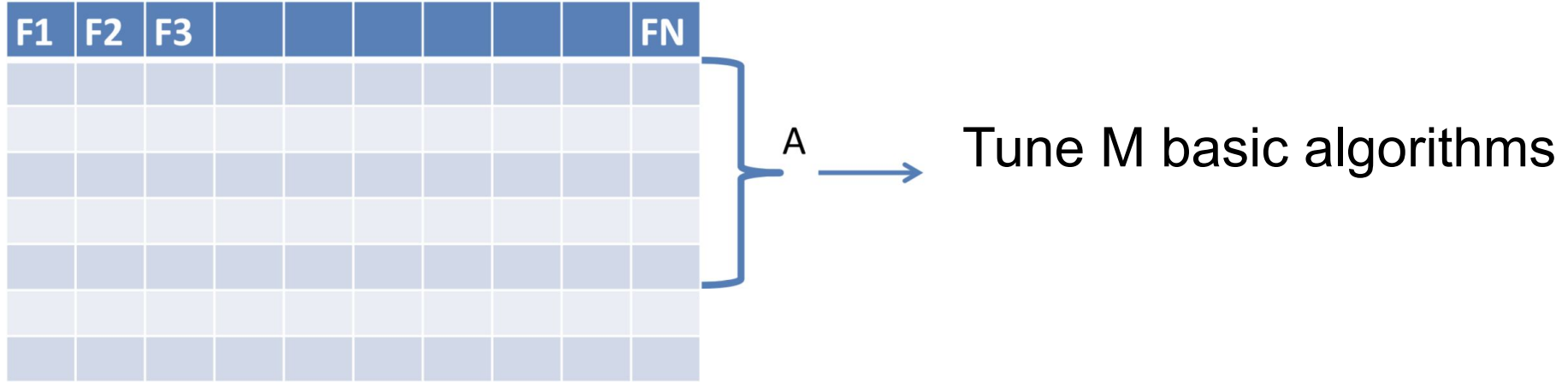
Never use `train_test_split` in this case!



# Stacking and blending

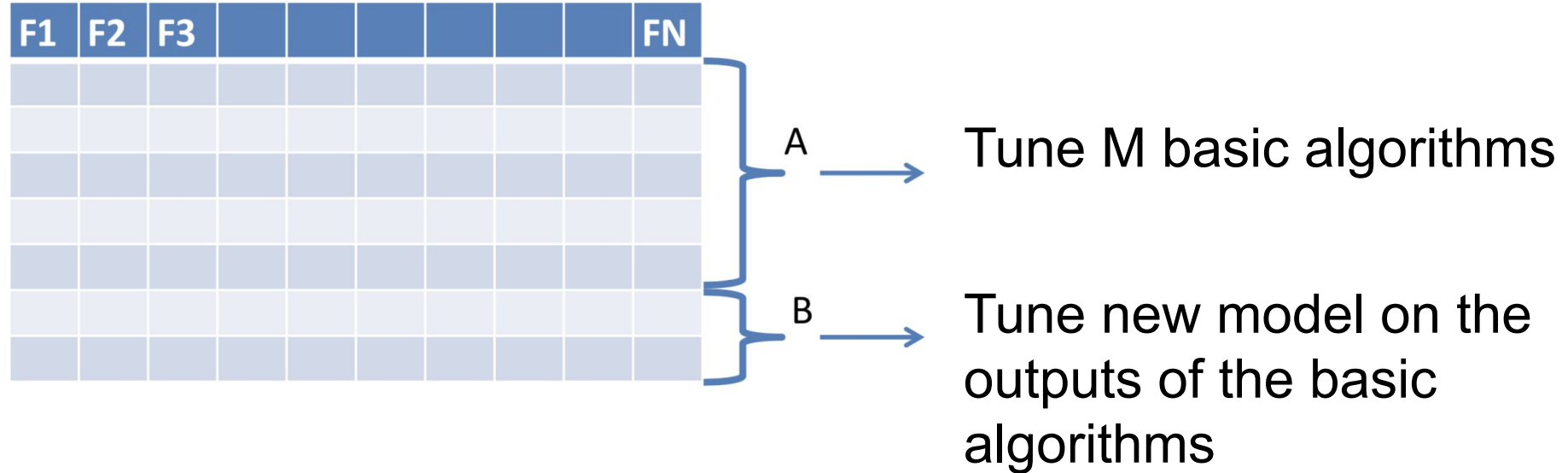
# Blending

How to build an ensemble from *different* models?



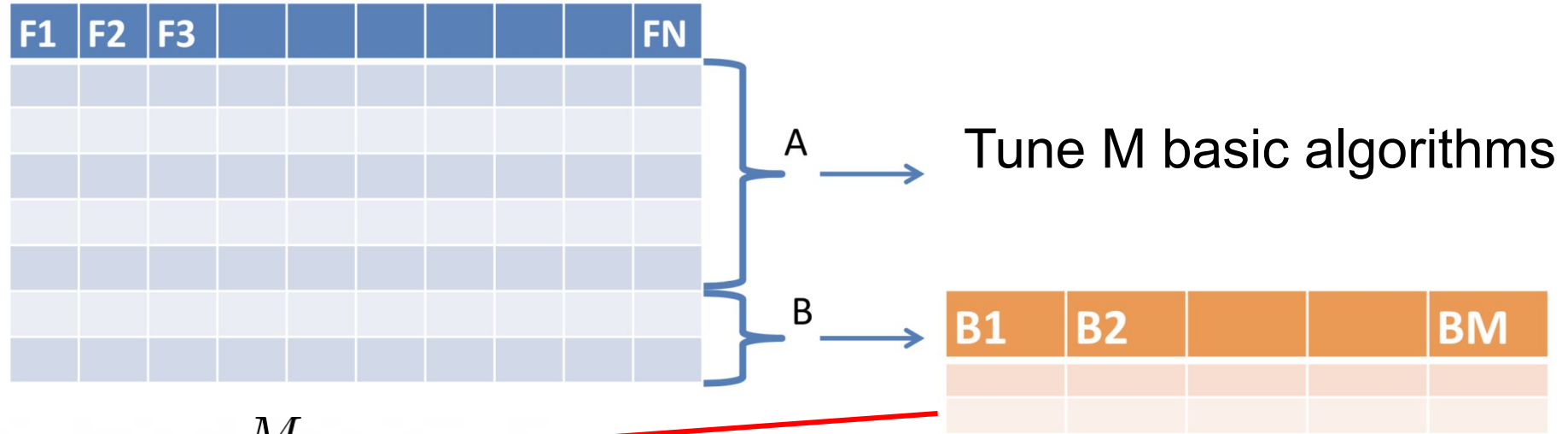
# Blending

How to build an ensemble from *different* models?



# Blending

How to build an ensemble from *different* models?



$$\hat{f}(x) = \sum_{i=1}^M \rho_i f_i(x)$$

$$\sum_{i=1}^M \rho_i = 1, \quad \rho_i \in [0; 1] \quad \forall i$$

# Blending

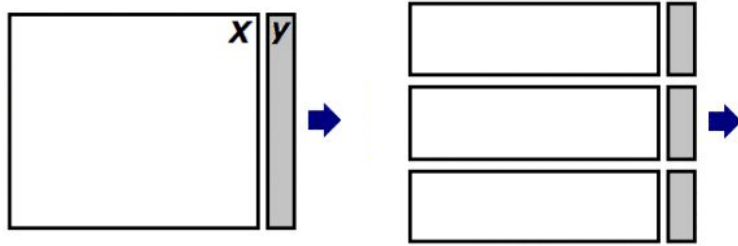
Just combine several *strong/complex* models.

$$\hat{f}(x) = \sum_{i=1}^M \rho_i f_i(x), \quad \sum_{i=1}^M \rho_i = 1, \quad \rho_i \in [0; 1] \quad \forall i$$

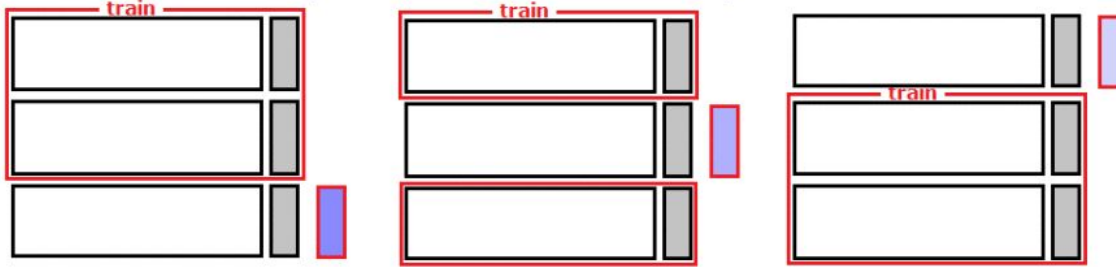
- Pros:
  - Simple and intuitive ensembling method.
  - Average several blendings to achieve better results.
- Cons:
  - Linear composition is not always enough.
  - Need to split the data. **How to fix it?**

# Stacking

## 1. Split data into folds



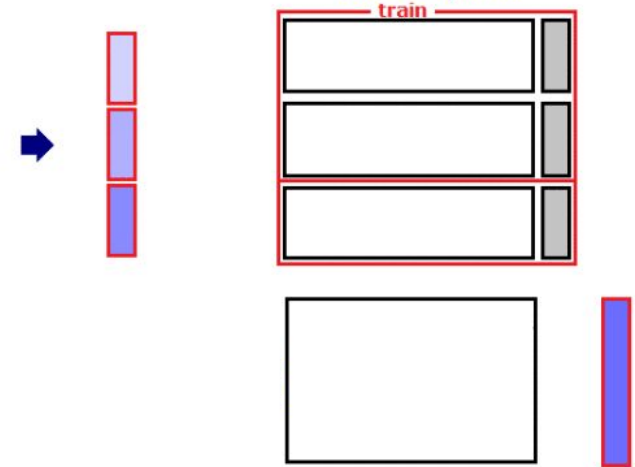
Fit using folds to get meta-features on train



## 2. Tune models on different groups of folds, predict on left out

## 3. Tune the new model on the “meta”-features

Fit using all data meta-features on test



# Stacking

- Train base algorithm(s) on different groups of folds leaving one fold out.
- Predict the meta-features on the left-out fold and test data.
- Train the meta-algorithm on the meta-features representation of the train data.
- Use it on the meta-features representation of the test data.

# Stacking

- Besides  $f_i(x)$ ,  $i=1,\dots,M$ ;  $\hat{f}(x)$  may also depend on:
  - original features
  - Internal representations in  $f_i(x)$  (e.g. class scores)



# Stacking

- Besides  $f_i(x)$ ,  $i=1,\dots,M$ ;  $\hat{f}(x)$  may also depend on:
  - original features
  - Internal representations in  $f_i(x)$  (e.g. class scores)

# Stacking

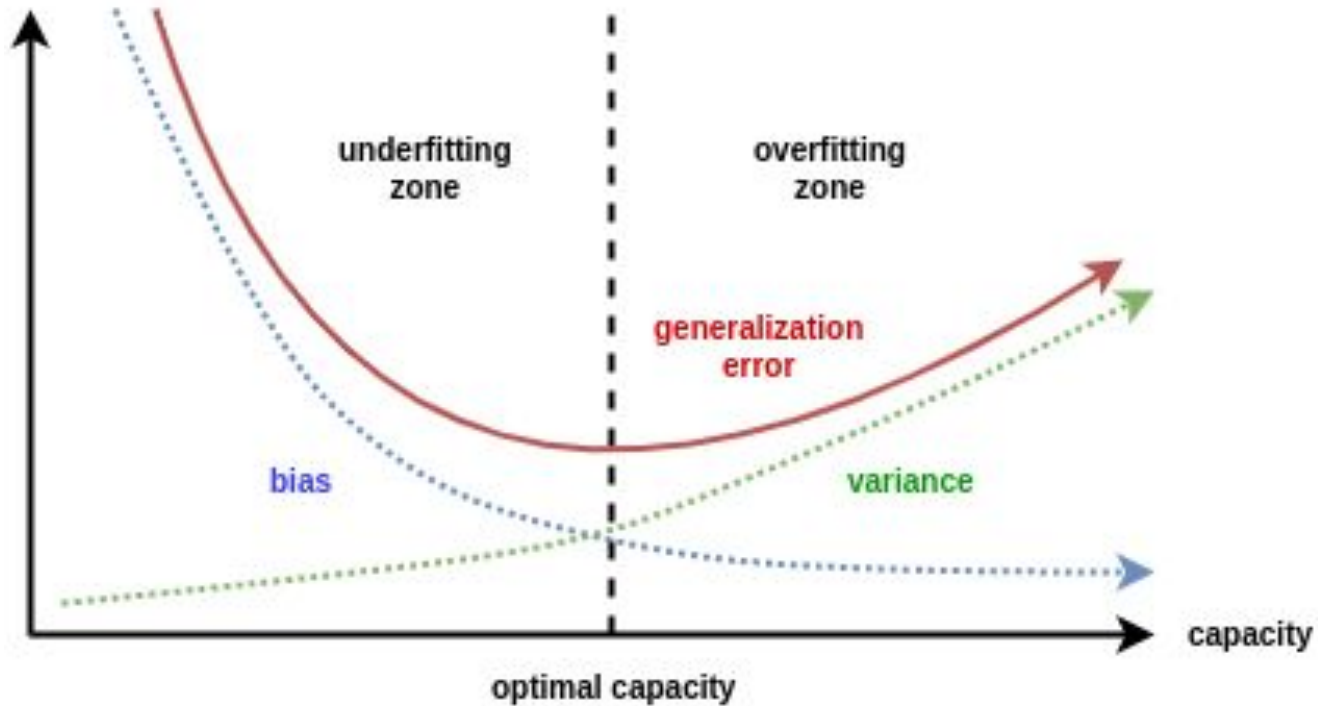
- Pros:
  - Powerful ensembling method, if you know how to use it
  - Quite popular in ML-competitions
  - One might perform stacking on the meta-features dataset as well
- Cons:
  - Meta-features on each fold are actually predicted by different models
    - However, regularization usually helps
  - Hard to explain your model behaviour

# Recap: ensembling methods as of now

1. Bagging.
2. Random subspace method (RSM).
3. Bagging + RSM + Decision trees = Random Forest.
4. Blending.
5. Stacking.

# Bias-Variance tradeoff

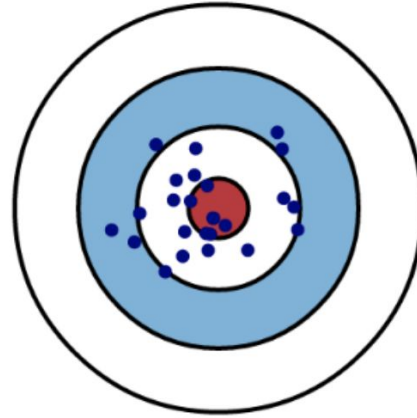
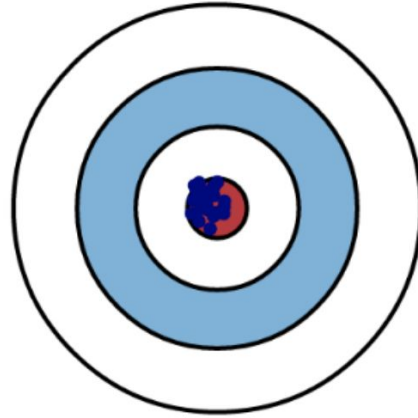
# Bias-variance tradeoff



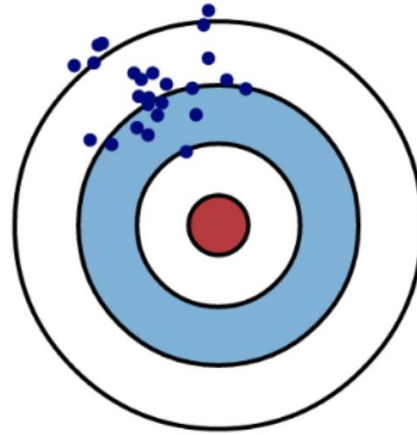
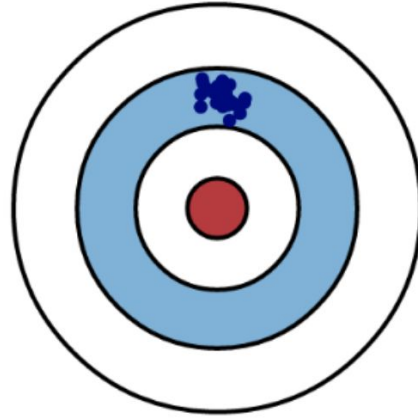
Low Variance

High Variance

Low Bias



High Bias



# Bias-variance decomposition derivation

# Bias-variance decomposition

The dataset  $X = (x_i, y_i)_{i=1}^{\ell}$  with  $y_i \in \mathbb{R}$   
for regression problem.

Denote loss function  $L(y, a) = (y - a(x))^2$ .

The empirical risk takes form:

$$R(a) = \mathbb{E}_{x,y} \left[ (y - a(x))^2 \right] = \int_{\mathbb{X}} \int_{\mathbb{Y}} p(x, y) (y - a(x))^2 dx dy.$$



# Bias-variance decomposition

Let's show that

$$a_*(x) = \mathbb{E}[y \mid x] = \int_{\mathbb{Y}} yp(y \mid x)dy = \arg \min_a R(a).$$

$$\begin{aligned} L(y, a(x)) &= (y - a(x))^2 = (y - \mathbb{E}(y \mid x) + \mathbb{E}(y \mid x) - a(x))^2 = \\ &= (y - \mathbb{E}(y \mid x))^2 + 2(y - \mathbb{E}(y \mid x))(\mathbb{E}(y \mid x) - a(x)) + (\mathbb{E}(y \mid x) - a(x))^2. \end{aligned}$$

Returning to the risk estimation:

$$\begin{aligned} R(a) &= \mathbb{E}_{x,y} L(y, a(x)) = \\ &= \mathbb{E}_{x,y} (y - \mathbb{E}(y \mid x))^2 + \mathbb{E}_{x,y} (\mathbb{E}(y \mid x) - a(x))^2 + \\ &+ 2\mathbb{E}_{x,y} (y - \mathbb{E}(y \mid x))(\mathbb{E}(y \mid x) - a(x)). \end{aligned}$$


$$R(a) = \mathbb{E}_{x,y} L(y, a(x)) =$$

Focus on the last term:

$$= \mathbb{E}_{x,y} (y - \mathbb{E}(y | x))^2 + \mathbb{E}_{x,y} (\mathbb{E}(y | x) - a(x))^2 +$$

$$+ 2\mathbb{E}_{x,y} (y - \mathbb{E}(y | x)) (\mathbb{E}(y | x) - a(x)).$$

Does not depend on y



$$\mathbb{E}_x \mathbb{E}_y \left[ (y - \mathbb{E}(y | x)) (\mathbb{E}(y | x) - a(x)) \mid x \right] =$$

$$= \mathbb{E}_x \left( (\mathbb{E}(y | x) - a(x)) \mathbb{E}_y \left[ (y - \mathbb{E}(y | x)) \mid x \right] \right) =$$

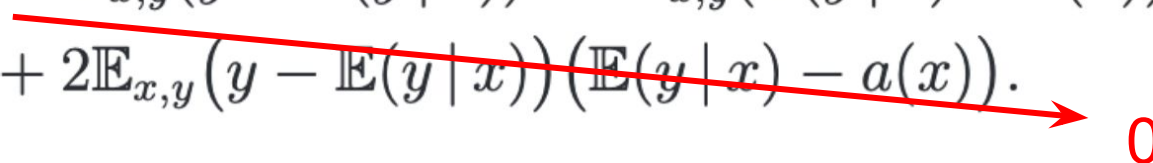
$$= \mathbb{E}_x \left( (\mathbb{E}(y | x) - a(x)) (\mathbb{E}(y | x) - \mathbb{E}(y | x)) \right) =$$

$$= 0$$

$$R(a) = \mathbb{E}_{x,y} L(y, a(x)) =$$

Focus on the last term:

$$= \mathbb{E}_{x,y} (y - \mathbb{E}(y | x))^2 + \mathbb{E}_{x,y} (\mathbb{E}(y | x) - a(x))^2 +$$

$$+ 2\mathbb{E}_{x,y} (y - \mathbb{E}(y | x)) (\mathbb{E}(y | x) - a(x)).$$


$$\mathbb{E}_x \mathbb{E}_y \left[ (y - \mathbb{E}(y | x)) (\mathbb{E}(y | x) - a(x)) \mid x \right] =$$

$$= \mathbb{E}_x \left( (\mathbb{E}(y | x) - a(x)) \mathbb{E}_y \left[ (y - \mathbb{E}(y | x)) \mid x \right] \right) =$$

$$= \mathbb{E}_x \left( (\mathbb{E}(y | x) - a(x)) (\mathbb{E}(y | x) - \mathbb{E}(y | x)) \right) =$$

$$= 0$$

So the risk takes form:

$$R(a) = \mathbb{E}_{x,y}(y - \mathbb{E}(y | x))^2 + \mathbb{E}_{x,y}(\mathbb{E}(y | x) - a(x))^2.$$

Does not depend on  $a(x)$

The minimum is reached when  $a(x) = \mathbb{E}(y | x)$ .

So the optimal regression model with square loss is

$$a_*(x) = \mathbb{E}(y | x) = \int_{\mathbb{Y}} yp(y | x)dy.$$

Denote  $\mu : (\mathbb{X} \times \mathbb{Y})^\ell \rightarrow \mathcal{A}$ , where  $\mathcal{A}$  is some family of algorithms.

Denote  $\mu : (\mathbb{X} \times \mathbb{Y})^\ell \rightarrow \mathcal{A}$ , where  $\mathcal{A}$  is some family of algorithms.

So  $L(\mu) = \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ (y - \mu(X)(x))^2 \right] \right]$ , where  $X$  dataset.

**In further slides (x) is omitted!**

Denote  $\mu : (\mathbb{X} \times \mathbb{Y})^\ell \rightarrow \mathcal{A}$ , where  $\mathcal{A}$  is some family of algorithms.

So  $L(\mu) = \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ (y - \mu(X)(x))^2 \right] \right]$ , where  $X$  dataset.

**In further slides (x) is omitted!**

If  $X$  is fixed, then

$$\mathbb{E}_{x,y} \left[ (y - \mu(X))^2 \right] = \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right] + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right].$$

Denote  $\mu : (\mathbb{X} \times \mathbb{Y})^\ell \rightarrow \mathcal{A}$ , where  $\mathcal{A}$  is some family of algorithms.

So  $L(\mu) = \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ (y - \mu(X)(x))^2 \right] \right]$ , where  $X$  dataset.

**In further slides (x) is omitted!**

If  $X$  is fixed, then

$$\mathbb{E}_{x,y} \left[ (y - \mu(X))^2 \right] = \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right] + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right].$$

Let's combine the latter equations:



Denote  $\mu : (\mathbb{X} \times \mathbb{Y})^\ell \rightarrow \mathcal{A}$ , where  $\mathcal{A}$  is some family of algorithms.

So  $L(\mu) = \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ (y - \mu(X)(x))^2 \right] \right]$ , where  $X$  dataset.

**In further slides (x) is omitted!**

If  $X$  is fixed, then

$$\mathbb{E}_{x,y} \left[ (y - \mu(X))^2 \right] = \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right] + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right].$$

Let's combine the latter equations:

$$L(\mu) = \mathbb{E}_X \left[ \underbrace{\mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right]}_{\text{Does not depend on } X} + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right]$$

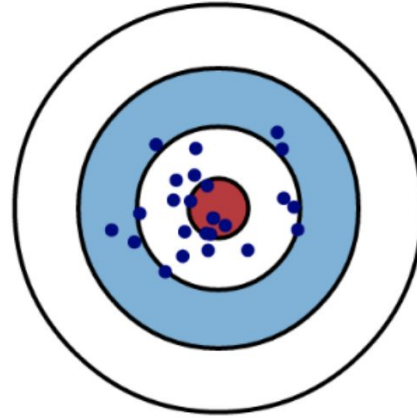
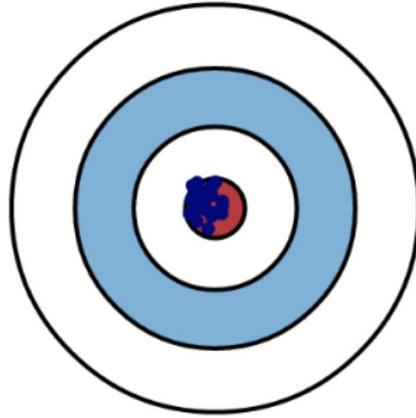
**Does not depend on  $X$**

$$\begin{aligned}
L(\mu) = & \underbrace{\mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right]}_{\text{noise}} + \\
& + \underbrace{\mathbb{E}_x \left[ (\mathbb{E}_X [\mu(X)] - \mathbb{E}[y | x])^2 \right]}_{\text{bias}} + \underbrace{\mathbb{E}_x \left[ \mathbb{E}_X \left[ (\mu(X) - \mathbb{E}_X [\mu(X)])^2 \right] \right]}_{\text{variance}}.
\end{aligned}$$

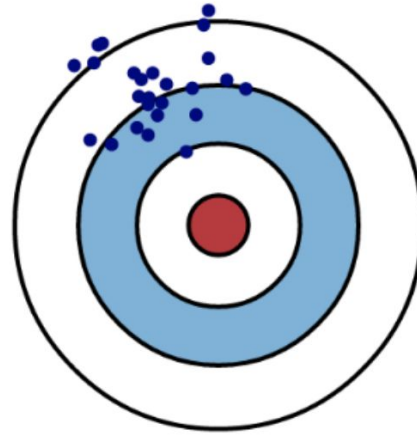
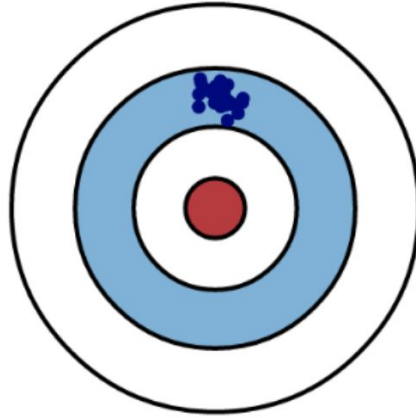
Low Variance

High Variance

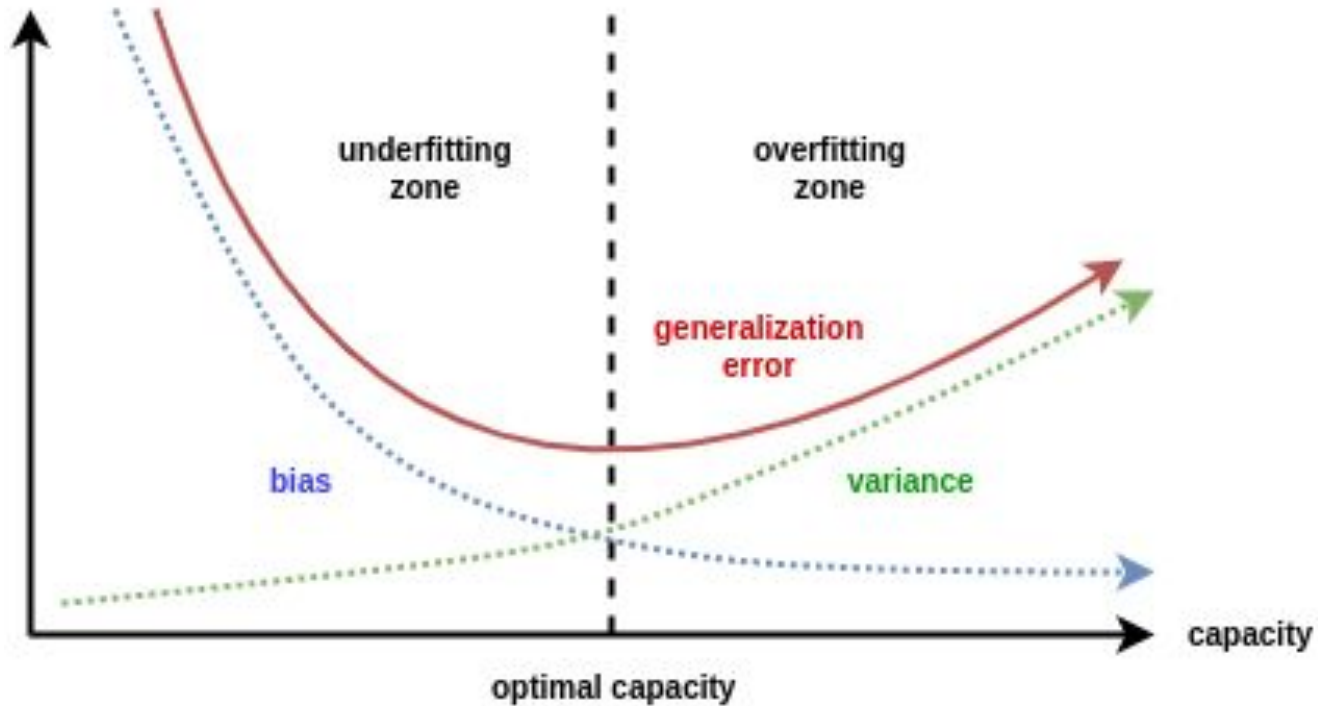
Low Bias



High Bias



# Bias-variance tradeoff



$$\begin{aligned}
 L(\mu) = & \underbrace{\mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right]}_{\text{noise}} + \\
 & \underbrace{\mathbb{E}_x \left[ (\mathbb{E}_X [\mu(X)] - \mathbb{E}[y | x])^2 \right]}_{\text{bias}} + \underbrace{\mathbb{E}_x \left[ \mathbb{E}_X \left[ (\mu(X) - \mathbb{E}_X [\mu(X)])^2 \right] \right]}_{\text{variance}}.
 \end{aligned}$$

This exact form of bias-variance decomposition is correct for square loss in regression.

However, it is much more general. See extra materials for more exotic cases.

Q & A



# Appendix

(continue from slide 17)



$$L(\mu) = \mathbb{E}_X \left[ \underbrace{\mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right]}_{\text{Does not depend on X}} + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right] =$$

Does not depend on X

$$L(\mu) = \mathbb{E}_X \left[ \underbrace{\mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right]}_{\text{Does not depend on X}} + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right] =$$

Does not depend on X

$$= \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right].$$

$$\begin{aligned} L(\mu) &= \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right] + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right] = \\ &= \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right]. \end{aligned}$$

Focus on the second term:

$$\begin{aligned}
 L(\mu) &= \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right] + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right] = \\
 &= \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right].
 \end{aligned}$$

Focus on the second term:

$$\mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right] =$$

$$\begin{aligned}
 L(\mu) &= \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right] + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right] = \\
 &= \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y | x])^2 \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right].
 \end{aligned}$$

Focus on the second term:

$$\begin{aligned}
 \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right] &= \\
 &= \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)] + \mathbb{E}_X [\mu(X)] - \mu(X))^2 \right] \right]
 \end{aligned}$$

$$\begin{aligned}\mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}[y \mid x] - \mu(X) \right)^2 \right] \right] &= \\ &= \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}[y \mid x] - \mathbb{E}_X [\mu(X)] + \mathbb{E}_X [\mu(X)] - \mu(X) \right)^2 \right] \right] =\end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}[y | x] - \mu(X) \right)^2 \right] \right] = \\
&= \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)] + \mathbb{E}_X [\mu(X)] - \mu(X) \right)^2 \right] \right] = \\
&= \mathbb{E}_{x,y} \left[ \underbrace{\mathbb{E}_X \left[ \left( \mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)] \right)^2 \right]} + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}_X [\mu(X)] - \mu(X) \right)^2 \right] \right] + \right. \\
&\quad \left. + 2 \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)] \right) \left( \mathbb{E}_X [\mu(X)] - \mu(X) \right) \right] \right] \right].
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right] = \\
&= \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)] + \mathbb{E}_X [\mu(X)] - \mu(X))^2 \right] \right] = \\
&= \mathbb{E}_{x,y} \left[ \underbrace{\mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)])^2 \right]}_{\text{Does not depend on X}} \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}_X [\mu(X)] - \mu(X))^2 \right] \right] + \\
&\quad + 2\mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)]) (\mathbb{E}_X [\mu(X)] - \mu(X)) \right] \right].
\end{aligned}$$



$$\begin{aligned}
& \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}[y | x] - \mu(X) \right)^2 \right] \right] = \\
&= \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)] + \mathbb{E}_X [\mu(X)] - \mu(X) \right)^2 \right] \right] = \\
&= \mathbb{E}_{x,y} \left[ \underbrace{\mathbb{E}_X \left[ \left( \mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)] \right)^2 \right]}_{\text{Does not depend on X}} \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}_X [\mu(X)] - \mu(X) \right)^2 \right] \right] + \\
&\quad + 2 \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)] \right) \left( \mathbb{E}_X [\mu(X)] - \mu(X) \right) \right] \right].
\end{aligned}$$

Just a bit further, we are almost there

$$\begin{aligned}
& \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right] = \\
&= \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)] + \mathbb{E}_X [\mu(X)] - \mu(X))^2 \right] \right] = \\
&= \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)])^2 \right] \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}_X [\mu(X)] - \mu(X))^2 \right] \right] + \\
&\quad + 2\mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mathbb{E}_X [\mu(X)]) (\mathbb{E}_X [\mu(X)] - \mu(X)) \right] \right].
\end{aligned}$$

Focus on this term

$$\mathbb{E}_X \left[ \left( \mathbb{E}[y \mid x] - \mathbb{E}_X [\mu(X)] \right) \left( \mathbb{E}_X [\mu(X)] - \mu(X) \right) \right] =$$

$$\begin{aligned}\mathbb{E}_X \left[ \left( \mathbb{E}[y \mid x] - \mathbb{E}_X [\mu(X)] \right) \left( \mathbb{E}_X [\mu(X)] - \mu(X) \right) \right] &= \\ &= \left( \mathbb{E}[y \mid x] - \mathbb{E}_X [\mu(X)] \right) \mathbb{E}_X \left[ \mathbb{E}_X [\mu(X)] - \mu(X) \right] =\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_X \left[ (\mathbb{E}[y \mid x] - \mathbb{E}_X [\mu(X)]) (\mathbb{E}_X [\mu(X)] - \mu(X)) \right] &= \\
&= (\mathbb{E}[y \mid x] - \mathbb{E}_X [\mu(X)]) \mathbb{E}_X \left[ \mathbb{E}_X [\mu(X)] - \mu(X) \right] = \\
&= (\mathbb{E}[y \mid x] - \mathbb{E}_X [\mu(X)]) \left[ \mathbb{E}_X [\mu(X)] - \mathbb{E}_X [\mu(X)] \right] =
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_X \left[ (\mathbb{E}[y \mid x] - \mathbb{E}_X [\mu(X)]) (\mathbb{E}_X [\mu(X)] - \mu(X)) \right] &= \\
&= (\mathbb{E}[y \mid x] - \mathbb{E}_X [\mu(X)]) \mathbb{E}_X [\mathbb{E}_X [\mu(X)] - \mu(X)] = \\
&= (\mathbb{E}[y \mid x] - \mathbb{E}_X [\mu(X)]) [\mathbb{E}_X [\mu(X)] - \mathbb{E}_X [\mu(X)]] = \\
&= 0.
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mu(X))^2 \right] \right] = \\
&= \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mathbb{E}_X[\mu(X)] + \mathbb{E}_X[\mu(X)] - \mu(X))^2 \right] \right] = \\
&= \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mathbb{E}_X[\mu(X)])^2 \right] \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}_X[\mu(X)] - \mu(X))^2 \right] \right] + \\
&\quad + 2\mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y | x] - \mathbb{E}_X[\mu(X)])(\mathbb{E}_X[\mu(X)] - \mu(X)) \right] \right].
\end{aligned}$$

0