

A crash course in machine learning

Florent Krzakala & Antoine Baker

https://sphinxteam.github.io/mlcourse_2019/

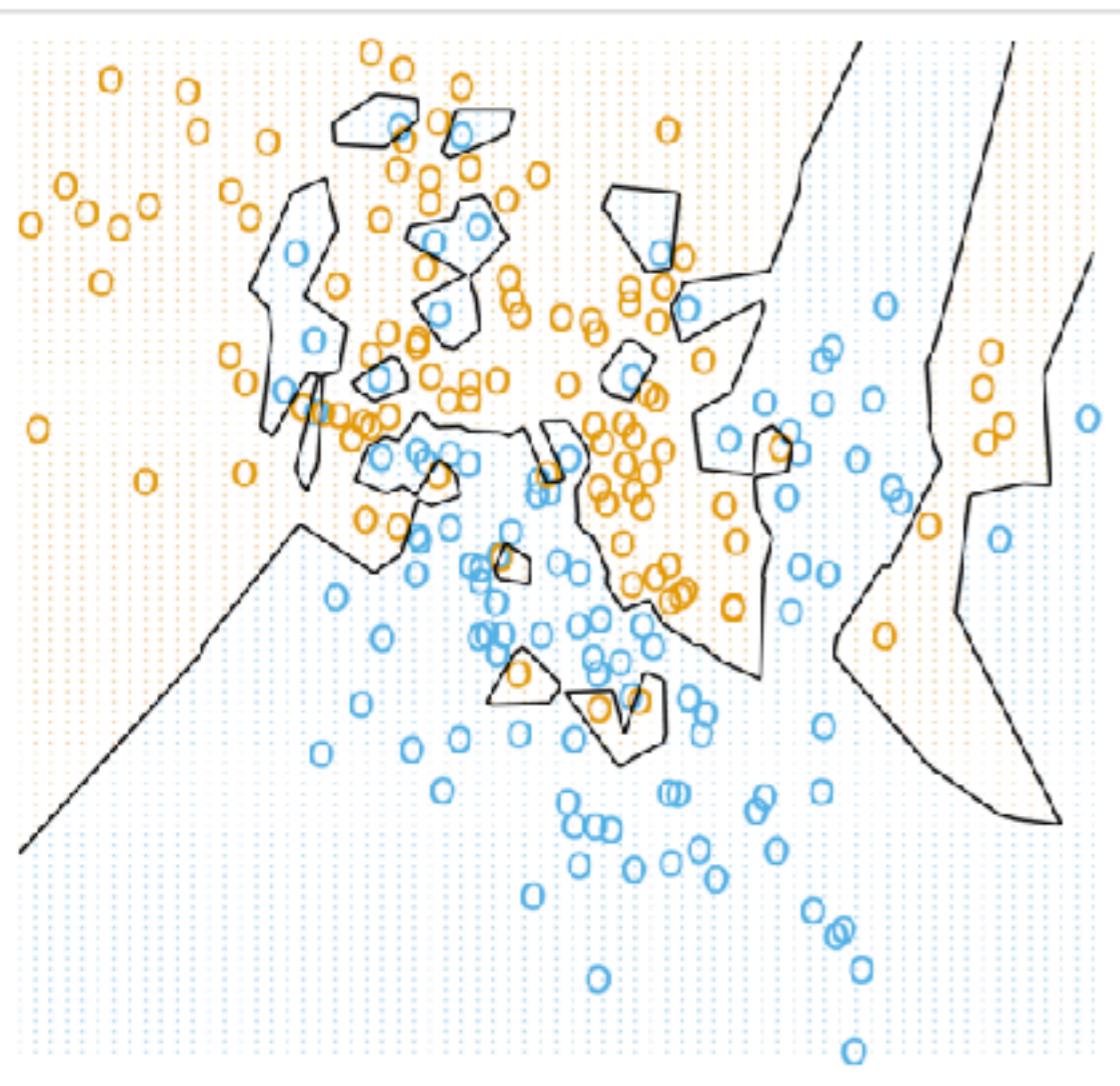
florent.krzakala@gmail.com

antoine.baker@gmail.com

So far...

Knn and linear methods

1-Nearest Neighbor Classifier



Solving $\textcolor{red}{y} \equiv \textcolor{blue}{A} \textcolor{green}{x} \in \mathbb{R}^m$ $\textcolor{blue}{A} \in \mathbb{R}^{m \times n}$

Determined ($m = n$): $\textcolor{blue}{x} = \textcolor{blue}{A}^{-1} \textcolor{red}{y}$

$$\textcolor{red}{y} = \begin{matrix} \textcolor{red}{y} \\ \vdots \\ \textcolor{red}{y} \end{matrix} = \begin{matrix} \textcolor{blue}{A} & \times & \textcolor{green}{x} \\ \hline \end{matrix}$$

Over-determined ($m > n$): $\min_{\textcolor{blue}{x}} \|\textcolor{blue}{A}\textcolor{blue}{x} - \textcolor{red}{y}\|^2$

$$\textcolor{blue}{x} = (\textcolor{blue}{A}^\top \textcolor{blue}{A})^{-1} \textcolor{blue}{A}^\top \textcolor{red}{y} \stackrel{\text{def.}}{=} \textcolor{blue}{A}^{-} \textcolor{red}{y}$$

$$\textcolor{red}{y} \approx \begin{matrix} \textcolor{red}{y} \\ \vdots \\ \textcolor{red}{y} \end{matrix} = \begin{matrix} \textcolor{blue}{A} & \times & \textcolor{green}{x} \\ \hline \end{matrix}$$

Under-determined ($m < n$): $\min_{\textcolor{blue}{x}} \{ \|\textcolor{blue}{x}\| ; \textcolor{blue}{A}\textcolor{blue}{x} = \textcolor{red}{y} \}$

$$\textcolor{blue}{x} = \textcolor{blue}{A}^\top (\textcolor{blue}{A}\textcolor{blue}{A}^\top)^{-1} \textcolor{red}{y} \stackrel{\text{def.}}{=} \textcolor{blue}{A}^{-} \textcolor{red}{y}$$

$$\textcolor{red}{y} = \begin{matrix} \textcolor{red}{y} \\ \vdots \\ \textcolor{red}{y} \end{matrix} = \begin{matrix} \textcolor{blue}{A} & \times & \textcolor{green}{x} \\ \hline \end{matrix}$$

A ill-posed and/or noise: $\min_{\textcolor{blue}{x}} \|\textcolor{blue}{A}\textcolor{blue}{x} - \textcolor{red}{y}\|^2 + \lambda \|\textcolor{blue}{x}\|^2$

$$\textcolor{blue}{x} = (\textcolor{blue}{A}^\top \textcolor{blue}{A} + \lambda \text{Id}_n)^{-1} \textcolor{blue}{A}^\top \textcolor{red}{y} \xrightarrow{\lambda \rightarrow 0} \textcolor{blue}{A}^+ \textcolor{red}{y}$$

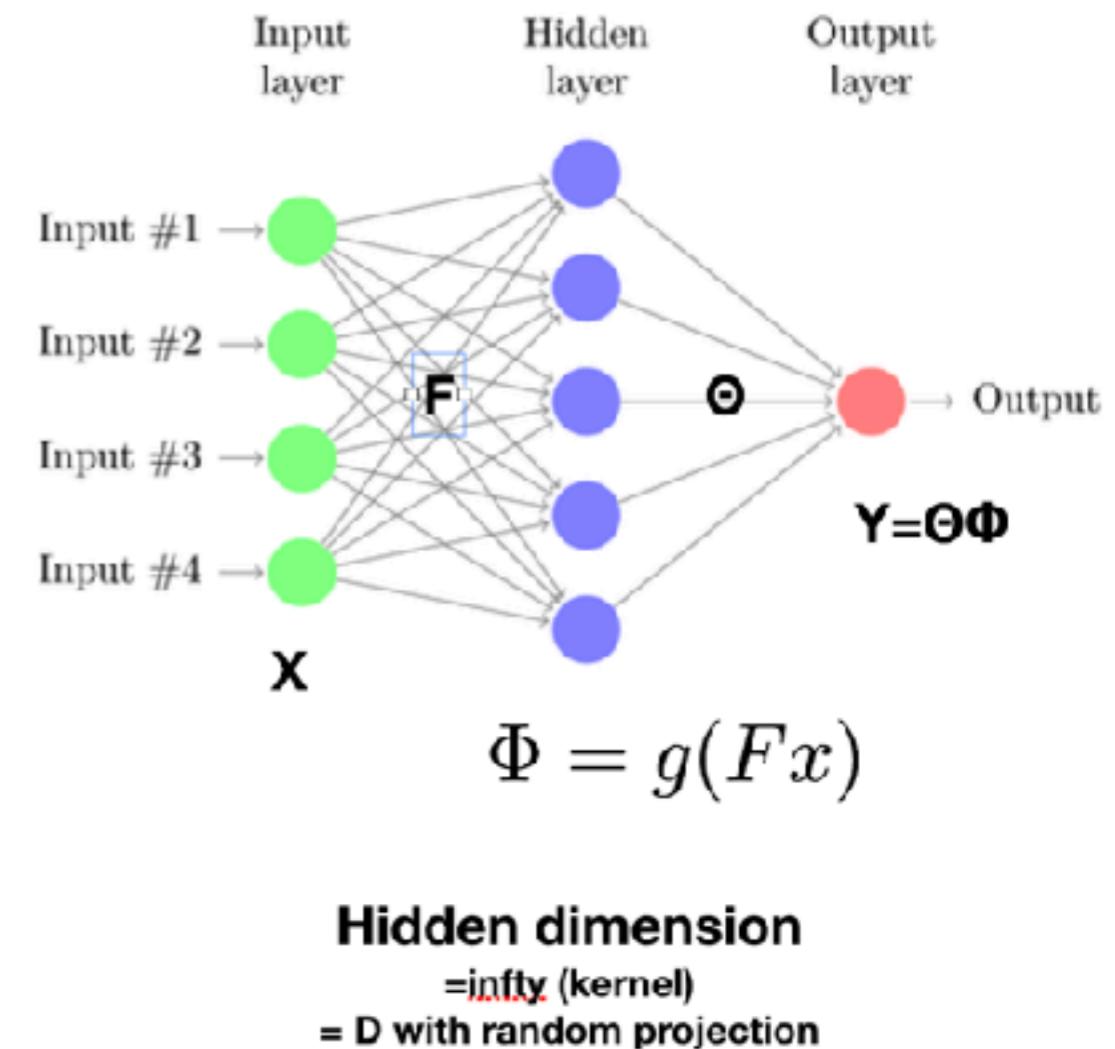
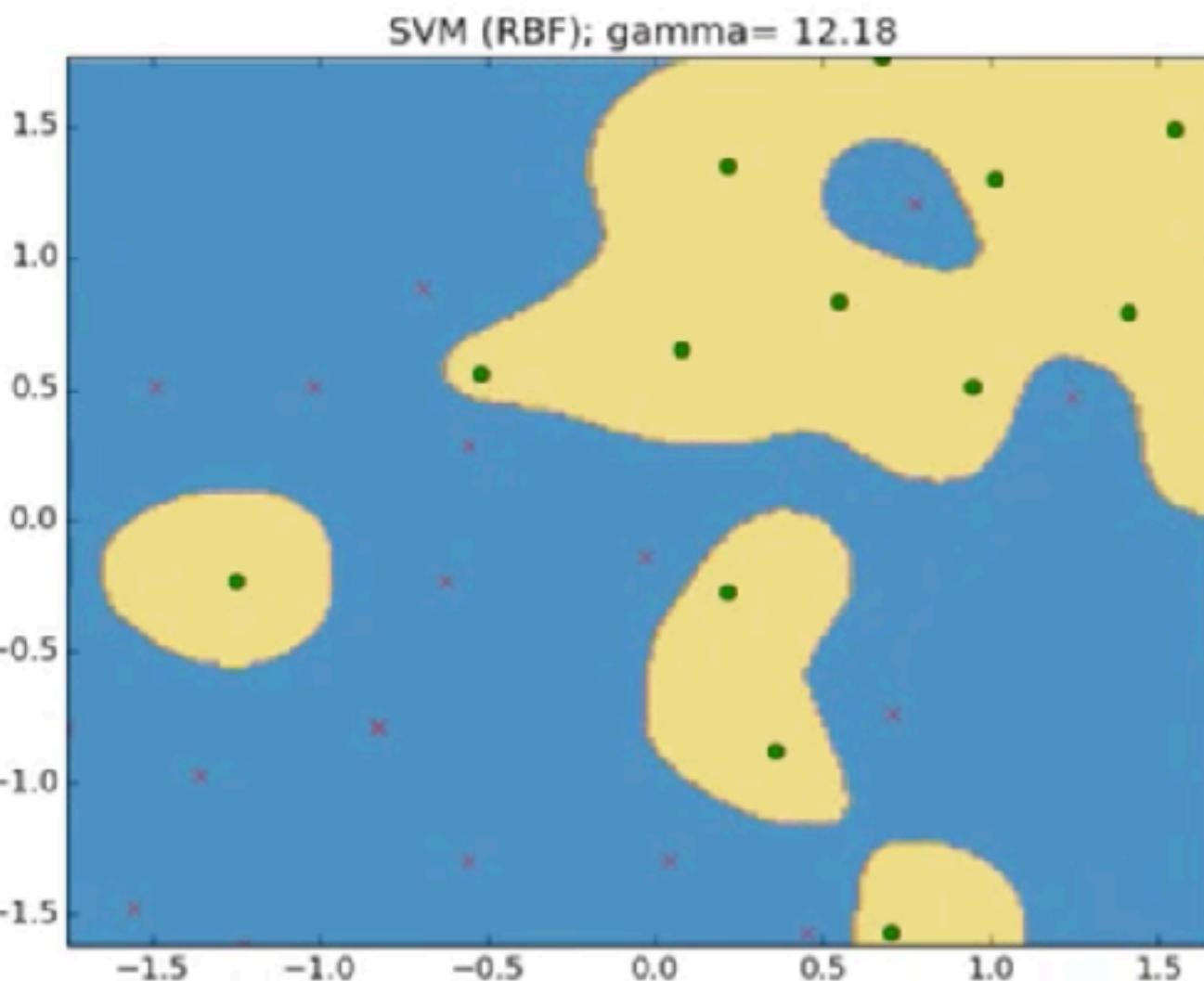
$$\textcolor{red}{y} \approx \begin{matrix} \textcolor{red}{y} \\ \vdots \\ \textcolor{red}{y} \end{matrix} = \begin{matrix} \textcolor{blue}{A} & \times & \textcolor{green}{x} \\ \hline \end{matrix}$$

$= \textcolor{blue}{A}^\top (\textcolor{blue}{A}\textcolor{blue}{A}^\top + \lambda \text{Id}_m)^{-1} \textcolor{red}{y}$ (Woodbury identity)

So far....

Kernels methods and random projections

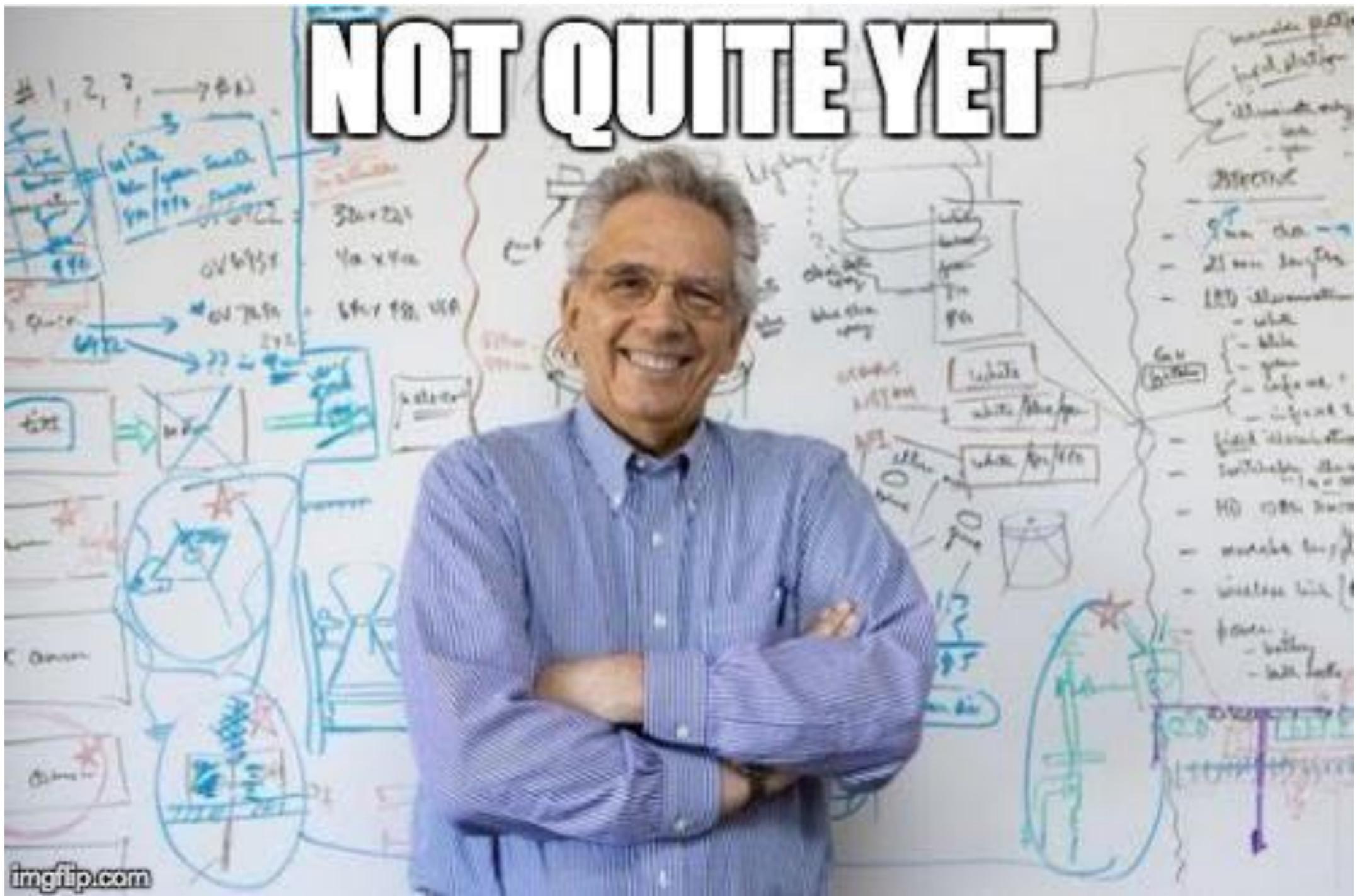
$$K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$





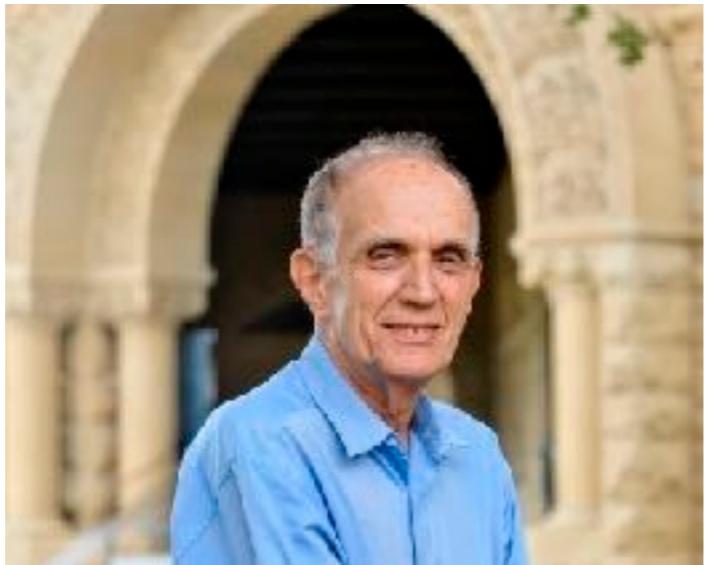
**That's it?
Are we finally going to speak about deep learning?**

NOT QUITE YET



Ensemble methods

Bagging (Bootstrap Aggregating)



Bradley Efron (1979)



Leo Breiman (1994)

Bootstrap

Why working with a single data set?

```
[ [ 0.524 ]  
[ 0.4 ]  
[ 0.361 ]  
[ 0.078 ]  
[ 0.859 ]  
[ 0.654 ]  
[ 0.384 ]  
[ 0.738 ]  
[ 0.156 ]  
[ 0.572 ] ]
```

Bootstrap

Why working with a single data set?

Create many data sets by sampling with replacement:

```
[ [ 0.524 ]  
[ 0.4 ]  
[ 0.361 ]  
[ 0.078 ]  
[ 0.859 ]  
[ 0.654 ]  
[ 0.384 ]  
[ 0.738 ]  
[ 0.156 ]  
[ 0.572 ] ]
```

Bootstrap

Why working with a single data set?

Create many data sets by sampling with replacement:

```
[ [ 0.524]  
[ 0.4 ]  
[ 0.361 ]  
[ 0.078 ]  
[ 0.859 ]  
[ 0.654 ]  
[ 0.384 ]  
[ 0.738 ]  
[ 0.156 ]  
[ 0.572 ] ]
```

```
[ [ 0.361 ]  
[ 0.572 ]  
[ 0.156 ]  
[ 0.361 ]  
[ 0.4 ]  
[ 0.384 ]  
[ 0.572 ]  
[ 0.654 ]  
[ 0.4 ]  
[ 0.384 ] ]
```

Bootstrap

Why working with a single data set?

Create many data sets by sampling with replacement:

```
[ [ 0.524]  
[ 0.4 ]  
[ 0.361]  
[ 0.078]  
[ 0.859]  
[ 0.654]  
[ 0.384]  
[ 0.738]  
[ 0.156]  
[ 0.572] ]
```

```
[ [ 0.361]  
[ 0.572]  
[ 0.156]  
[ 0.361]  
[ 0.4 ]  
[ 0.384]  
[ 0.572]  
[ 0.654]  
[ 0.4 ]  
[ 0.384 ] ]
```

```
[ [ 0.361]  
[ 0.156]  
[ 0.4 ]  
[ 0.654]  
[ 0.654 ]  
[ 0.859 ]  
[ 0.384 ]  
[ 0.4 ]  
[ 0.524 ]  
[ 0.859 ] ]
```

Bootstrap

Why working with a single data set?

Create many data sets by sampling with replacement:

```
[ [ 0.524 ]  
[ 0.4 ]  
[ 0.361 ]  
[ 0.078 ]  
[ 0.859 ]  
[ 0.654 ]  
[ 0.384 ]  
[ 0.738 ]  
[ 0.156 ]  
[ 0.572 ] ]
```

```
[ [ 0.361 ]  
[ 0.572 ]  
[ 0.156 ]  
[ 0.361 ]  
[ 0.4 ]  
[ 0.384 ]  
[ 0.572 ]  
[ 0.654 ]  
[ 0.4 ]  
[ 0.384 ] ]
```

```
[ [ 0.361 ]  
[ 0.156 ]  
[ 0.4 ]  
[ 0.654 ]  
[ 0.654 ]  
[ 0.859 ]  
[ 0.384 ]  
[ 0.4 ]  
[ 0.524 ]  
[ 0.859 ] ]
```

```
[ [ 0.156 ]  
[ 0.156 ]  
[ 0.524 ]  
[ 0.078 ]  
[ 0.524 ]  
[ 0.859 ]  
[ 0.361 ]  
[ 0.078 ]  
[ 0.654 ]  
[ 0.4 ] ]
```

Bootstrap

Why working with a single data set?

Create many data sets by sampling with replacement:

```
[ [ 0.524]  
[ 0.4 ]  
[ 0.361 ]  
[ 0.078 ]  
[ 0.859 ]  
[ 0.654 ]  
[ 0.384 ]  
[ 0.738 ]  
[ 0.156 ]  
[ 0.572 ] ]
```

```
[ [ 0.361 ]  
[ 0.572 ]  
[ 0.156 ]  
[ 0.361 ]  
[ 0.4 ]  
[ 0.384 ]  
[ 0.572 ]  
[ 0.654 ]  
[ 0.4 ]  
[ 0.384 ] ]
```

```
[ [ 0.361 ]  
[ 0.156 ]  
[ 0.4 ]  
[ 0.654 ]  
[ 0.654 ]  
[ 0.384 ]  
[ 0.4 ]  
[ 0.524 ]  
[ 0.859 ] ]
```

```
[ [ 0.156 ]  
[ 0.156 ]  
[ 0.524 ]  
[ 0.078 ]  
[ 0.524 ]  
[ 0.859 ]  
[ 0.361 ]  
[ 0.078 ]  
[ 0.654 ]  
[ 0.4 ] ]
```

```
[ [ 0.078 ]  
[ 0.524 ]  
[ 0.859 ]  
[ 0.384 ]  
[ 0.859 ]  
[ 0.361 ]  
[ 0.361 ]  
[ 0.738 ]  
[ 0.572 ]  
[ 0.156 ] ]
```

Bootstrap

Why working with a single data set?

Create many data sets by sampling with replacement:

```
[ [ 0.524]  
[ 0.4 ]  
[ 0.361 ]  
[ 0.078 ]  
[ 0.859 ]  
[ 0.654 ]  
[ 0.384 ]  
[ 0.738 ]  
[ 0.156 ]  
[ 0.572 ] ]
```

```
[ [ 0.361 ]  
[ 0.572 ]  
[ 0.156 ]  
[ 0.361 ]  
[ 0.4 ]  
[ 0.384 ]  
[ 0.572 ]  
[ 0.654 ]  
[ 0.4 ]  
[ 0.384 ] ]
```

```
[ [ 0.361 ]  
[ 0.156 ]  
[ 0.4 ]  
[ 0.654 ]  
[ 0.654 ]  
[ 0.384 ]  
[ 0.4 ]  
[ 0.524 ]  
[ 0.859 ] ]
```

```
[ [ 0.156 ]  
[ 0.156 ]  
[ 0.524 ]  
[ 0.078 ]  
[ 0.524 ]  
[ 0.859 ]  
[ 0.361 ]  
[ 0.654 ]  
[ 0.4 ] ]
```

```
[ [ 0.078 ]  
[ 0.524 ]  
[ 0.859 ]  
[ 0.384 ]  
[ 0.859 ]  
[ 0.361 ]  
[ 0.361 ]  
[ 0.572 ]  
[ 0.156 ] ]
```

Consider these datasets are equally valid

Bootstrap

Why working with a single data set?

Create many data sets by sampling with replacement:

```
[ [ 0.524]  
[ 0.4 ]  
[ 0.361 ]  
[ 0.078 ]  
[ 0.859 ]  
[ 0.654 ]  
[ 0.384 ]  
[ 0.738 ]  
[ 0.156 ]  
[ 0.572 ] ]
```

```
[ [ 0.361 ]  
[ 0.572 ]  
[ 0.156 ]  
[ 0.361 ]  
[ 0.4 ]  
[ 0.384 ]  
[ 0.572 ]  
[ 0.654 ]  
[ 0.4 ]  
[ 0.384 ] ]
```

```
[ [ 0.361 ]  
[ 0.156 ]  
[ 0.4 ]  
[ 0.654 ]  
[ 0.654 ]  
[ 0.384 ]  
[ 0.4 ]  
[ 0.524 ]  
[ 0.859 ] ]
```

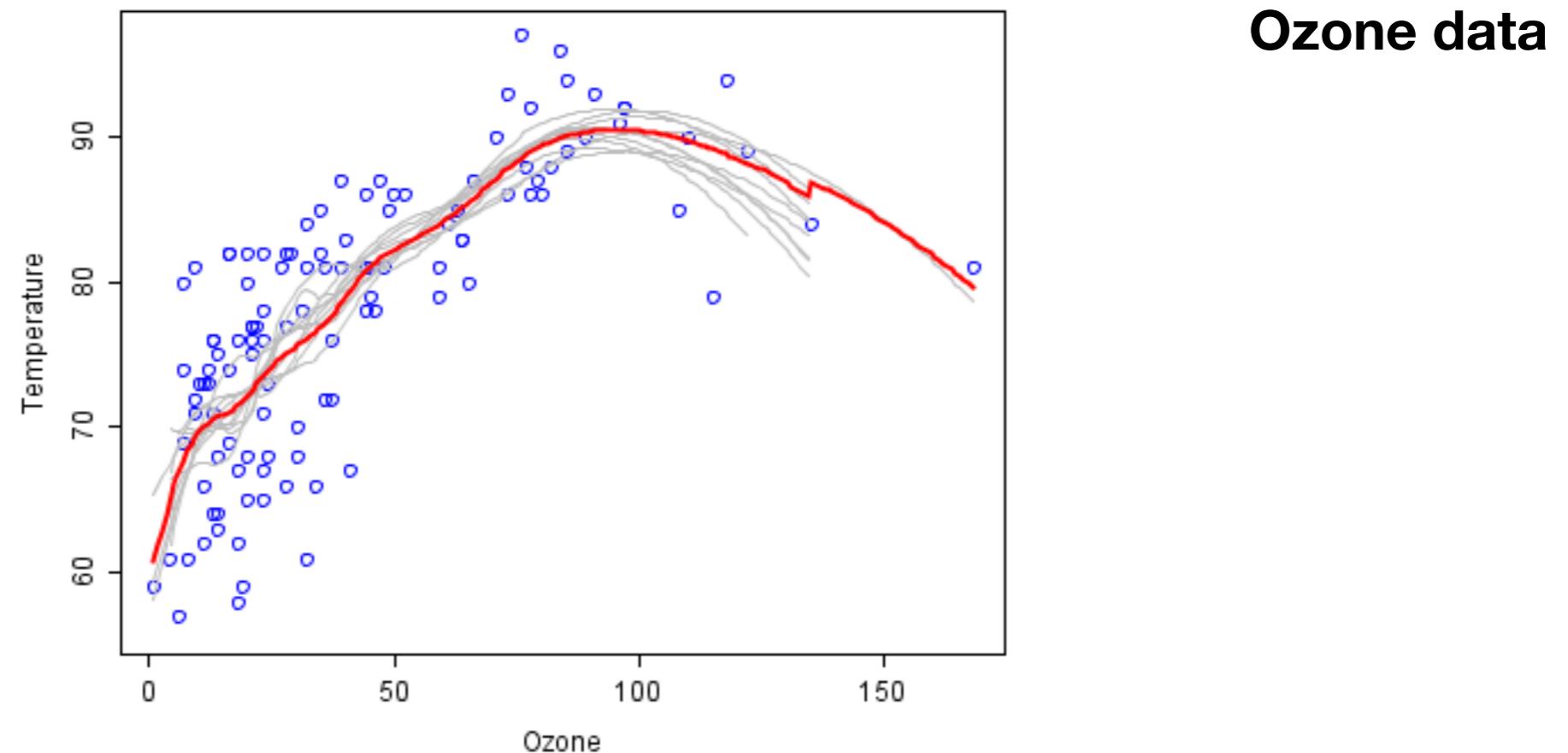
```
[ [ 0.156 ]  
[ 0.156 ]  
[ 0.524 ]  
[ 0.078 ]  
[ 0.524 ]  
[ 0.859 ]  
[ 0.361 ]  
[ 0.078 ]  
[ 0.654 ]  
[ 0.4 ] ]
```

```
[ [ 0.078 ]  
[ 0.524 ]  
[ 0.859 ]  
[ 0.384 ]  
[ 0.859 ]  
[ 0.361 ]  
[ 0.361 ]  
[ 0.738 ]  
[ 0.572 ]  
[ 0.156 ] ]
```

Consider these datasets are equally valid

Many uses: error analysis, statistical tests, ...

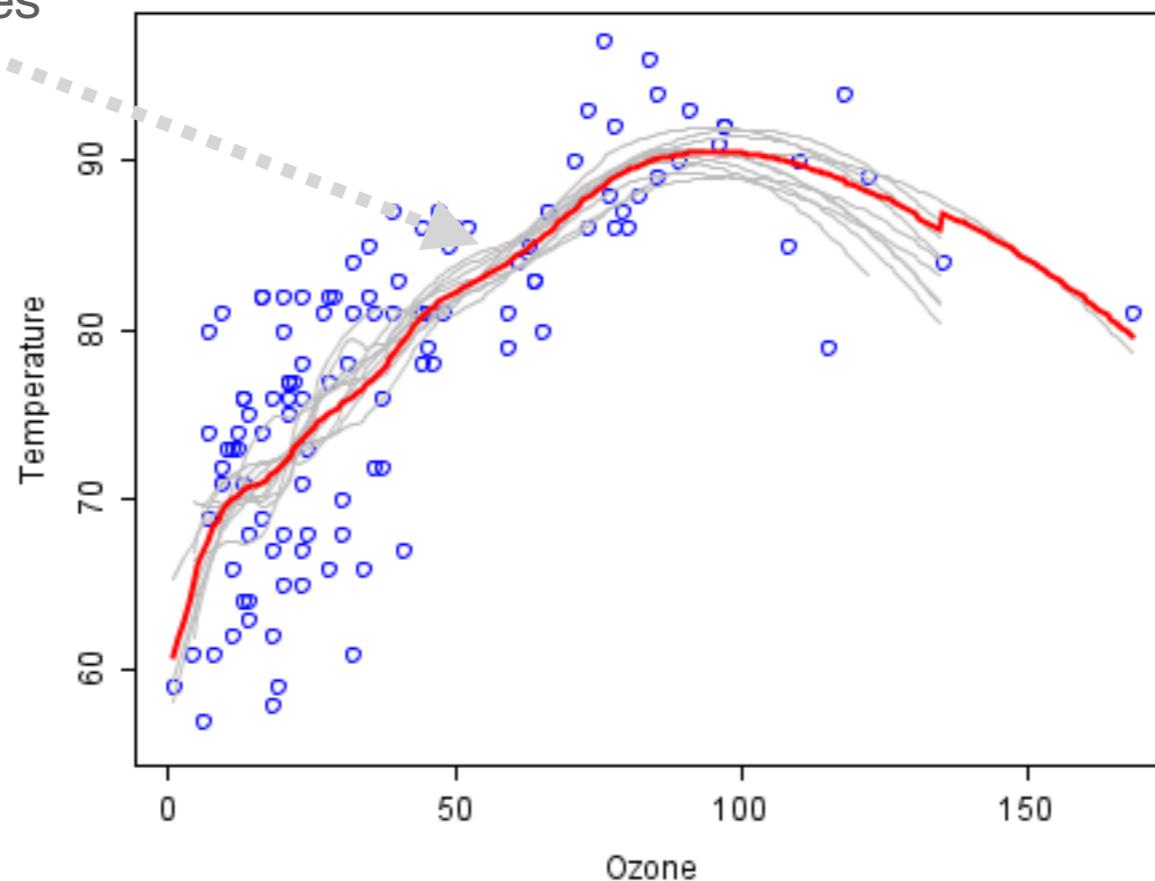
Bagging (Bootstrap Aggregating)



Bagging

(Bootstrap Aggregating)

Fits on bootstrap samples

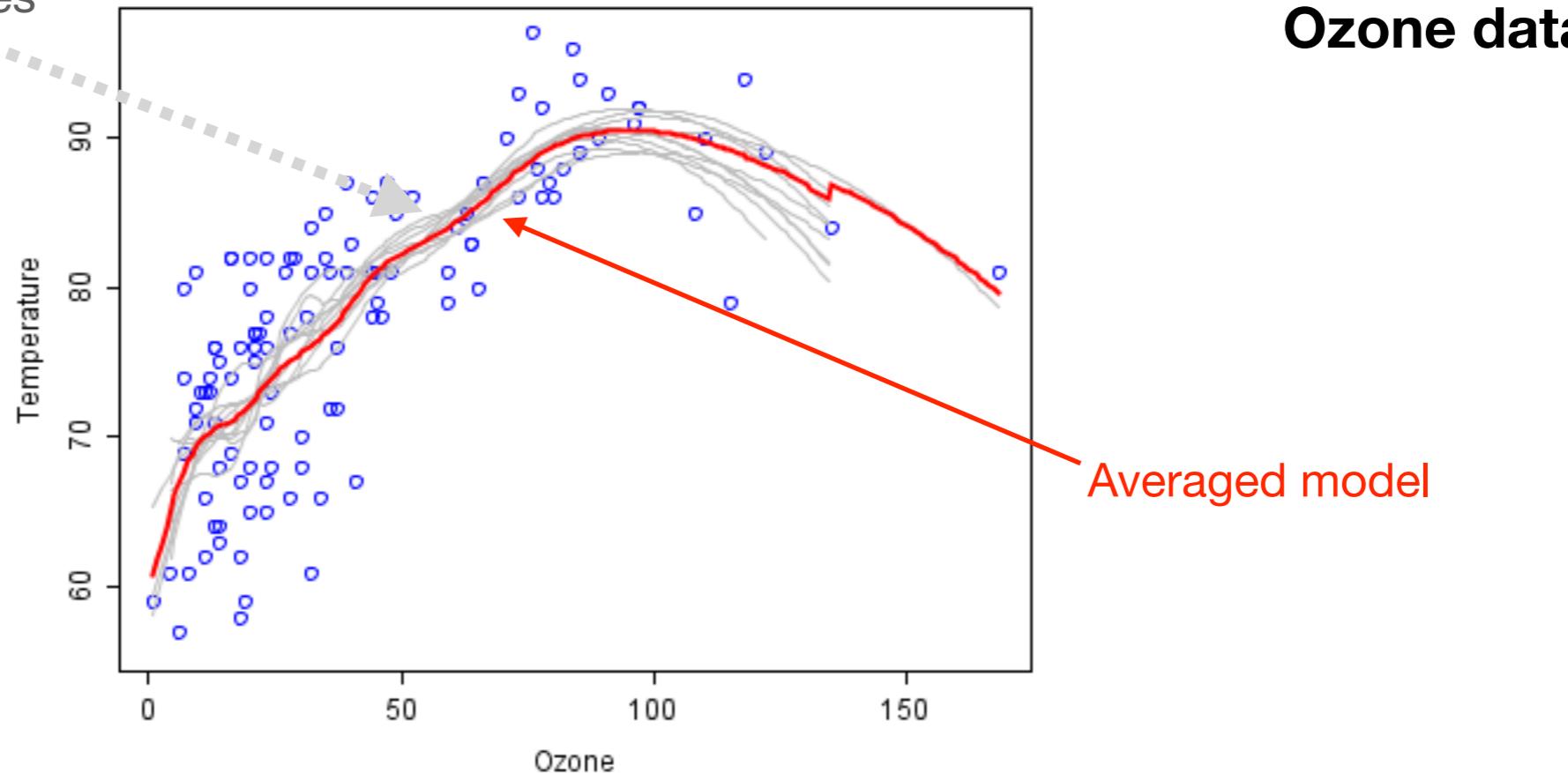


Ozone data

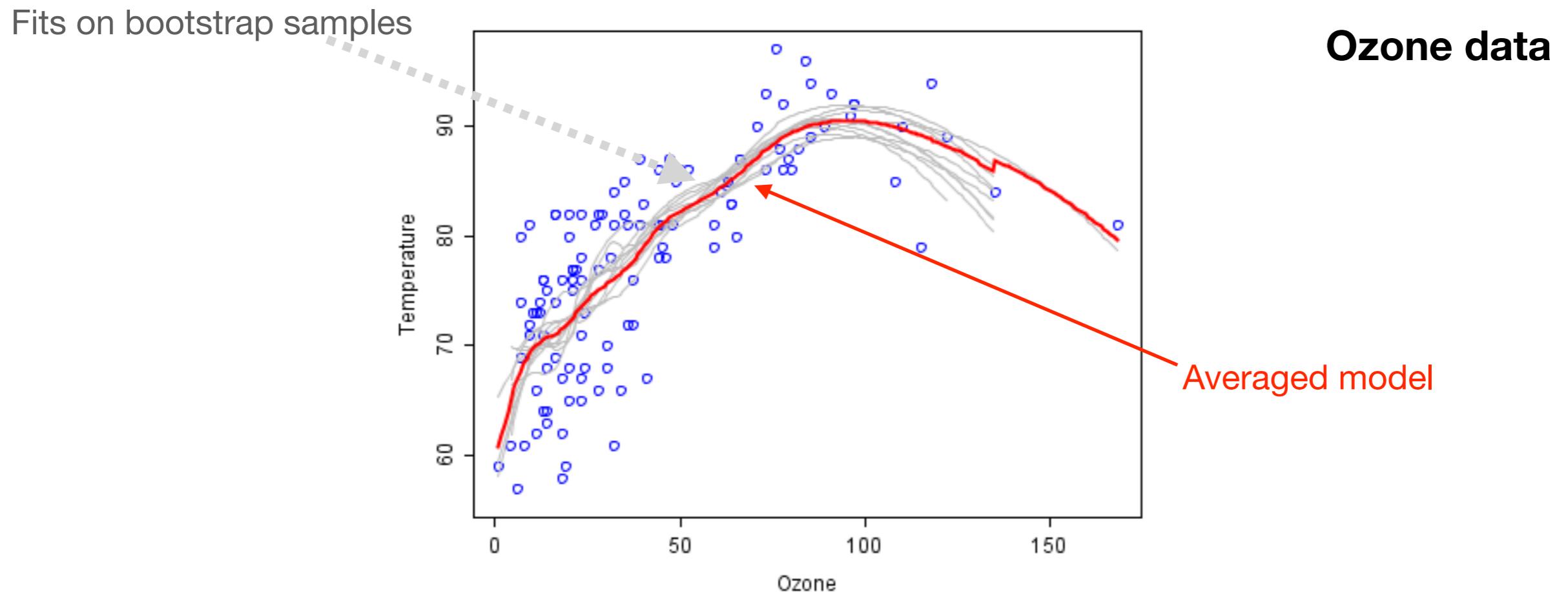
Bagging

(Bootstrap Aggregating)

Fits on bootstrap samples



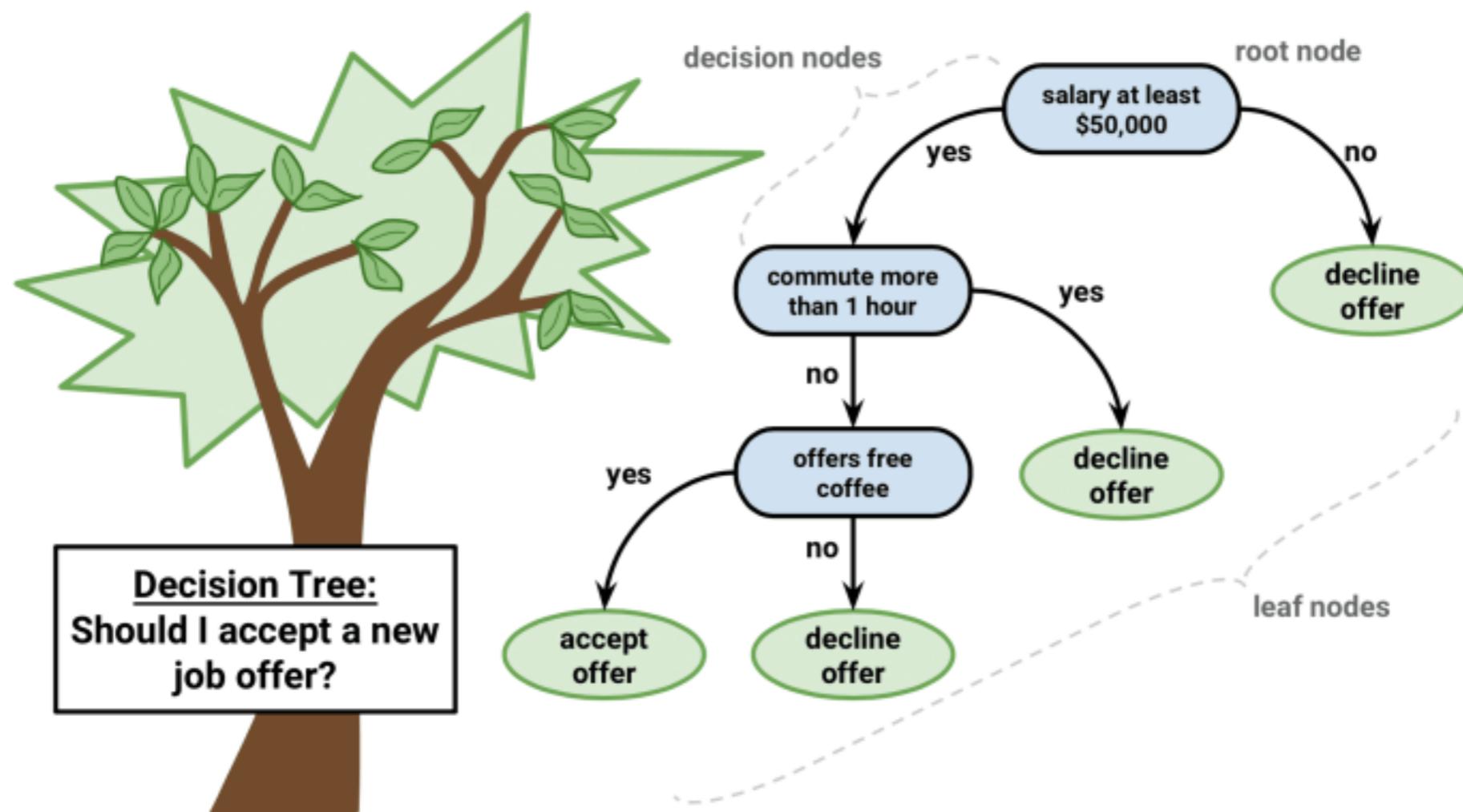
Bagging (Bootstrap Aggregating)

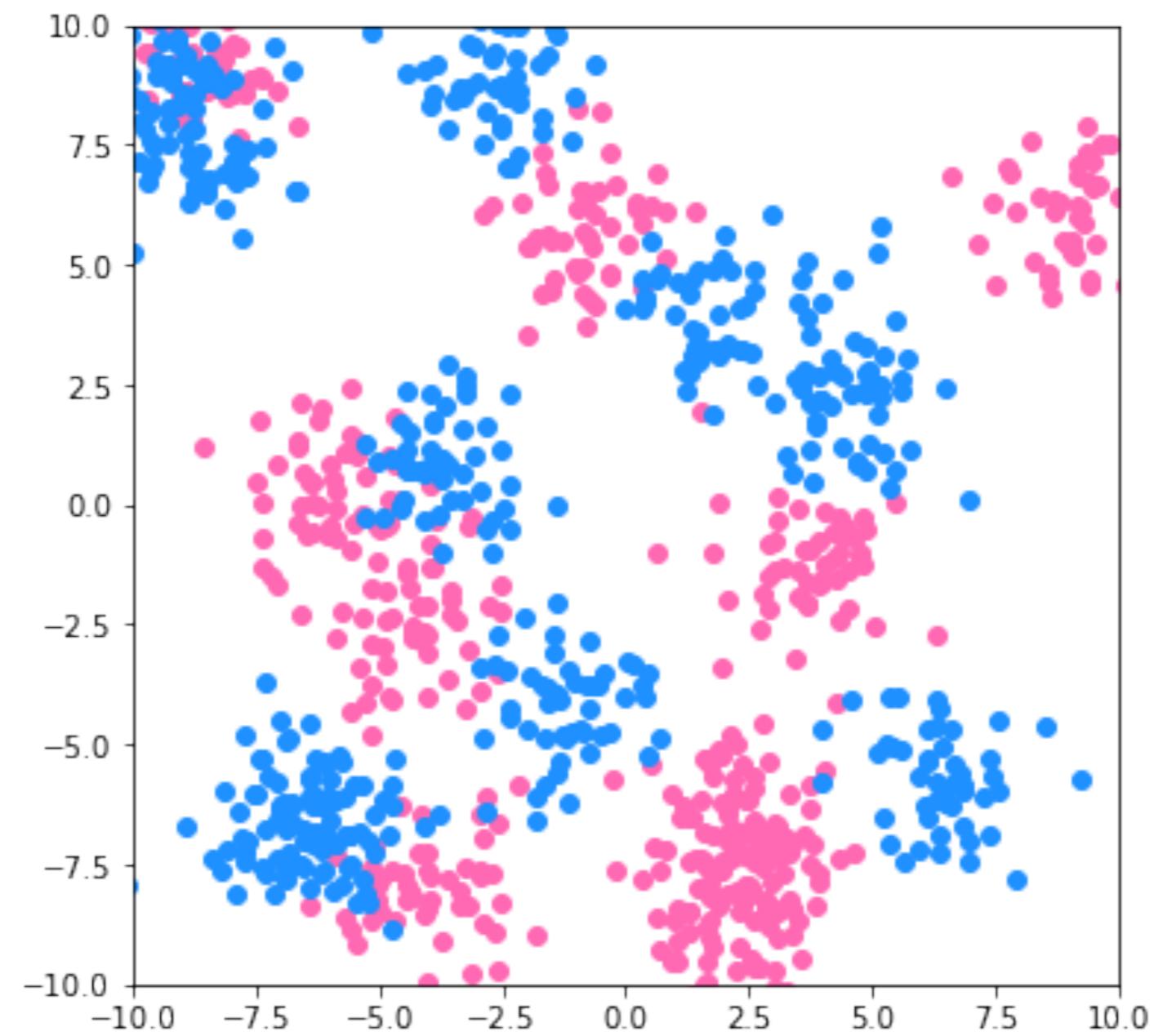


Averaging models reduces overfitting!

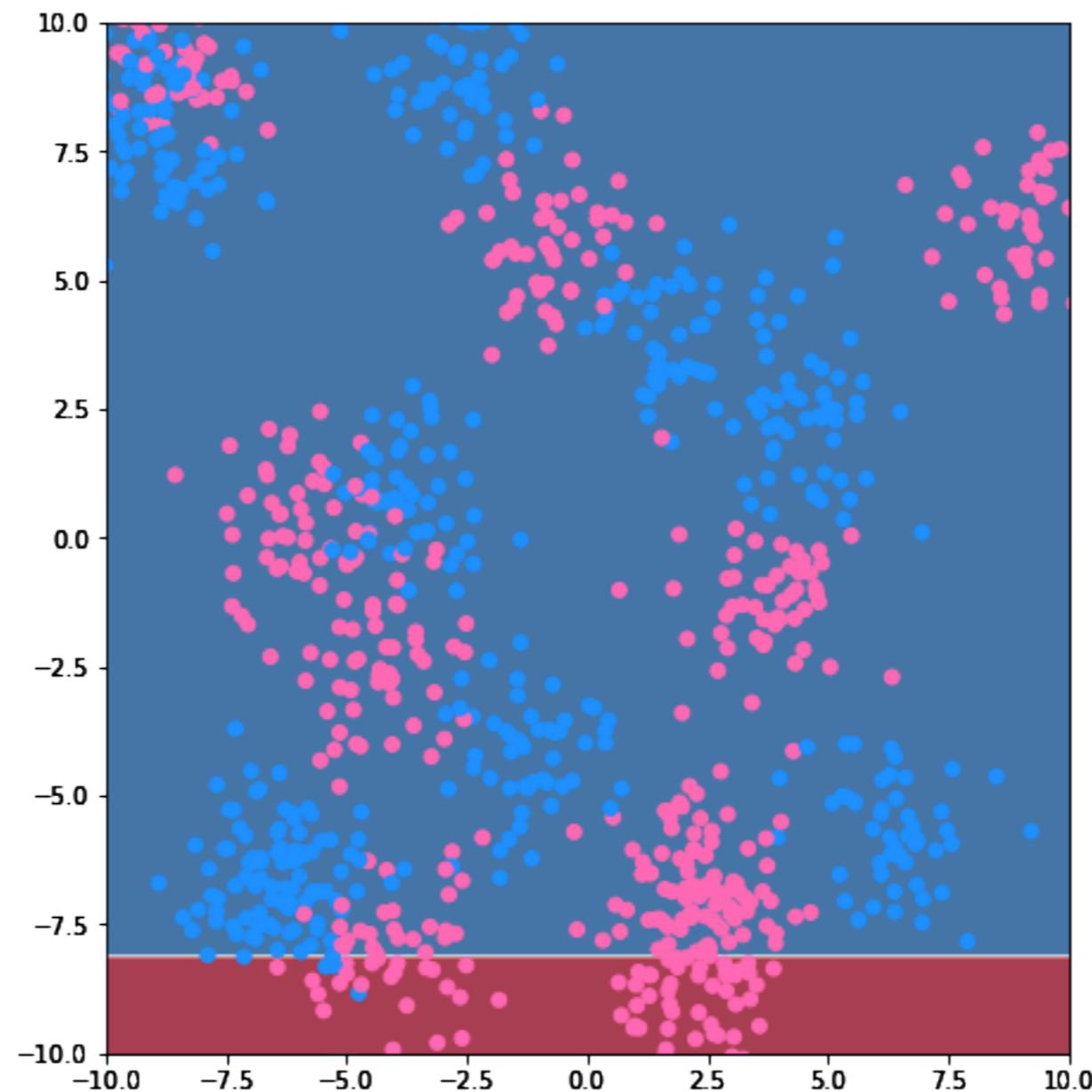
And for classification?

Decision trees!

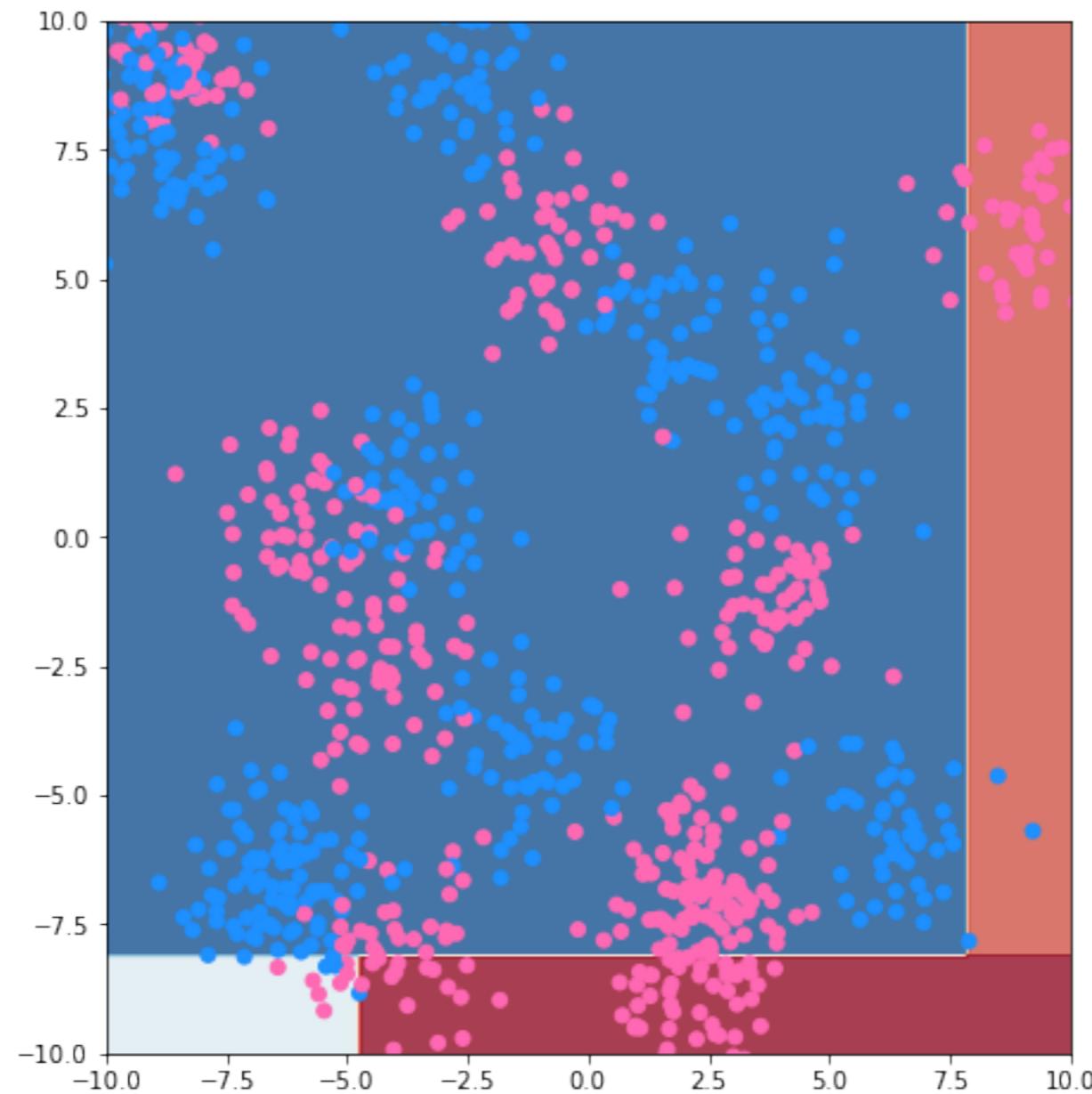




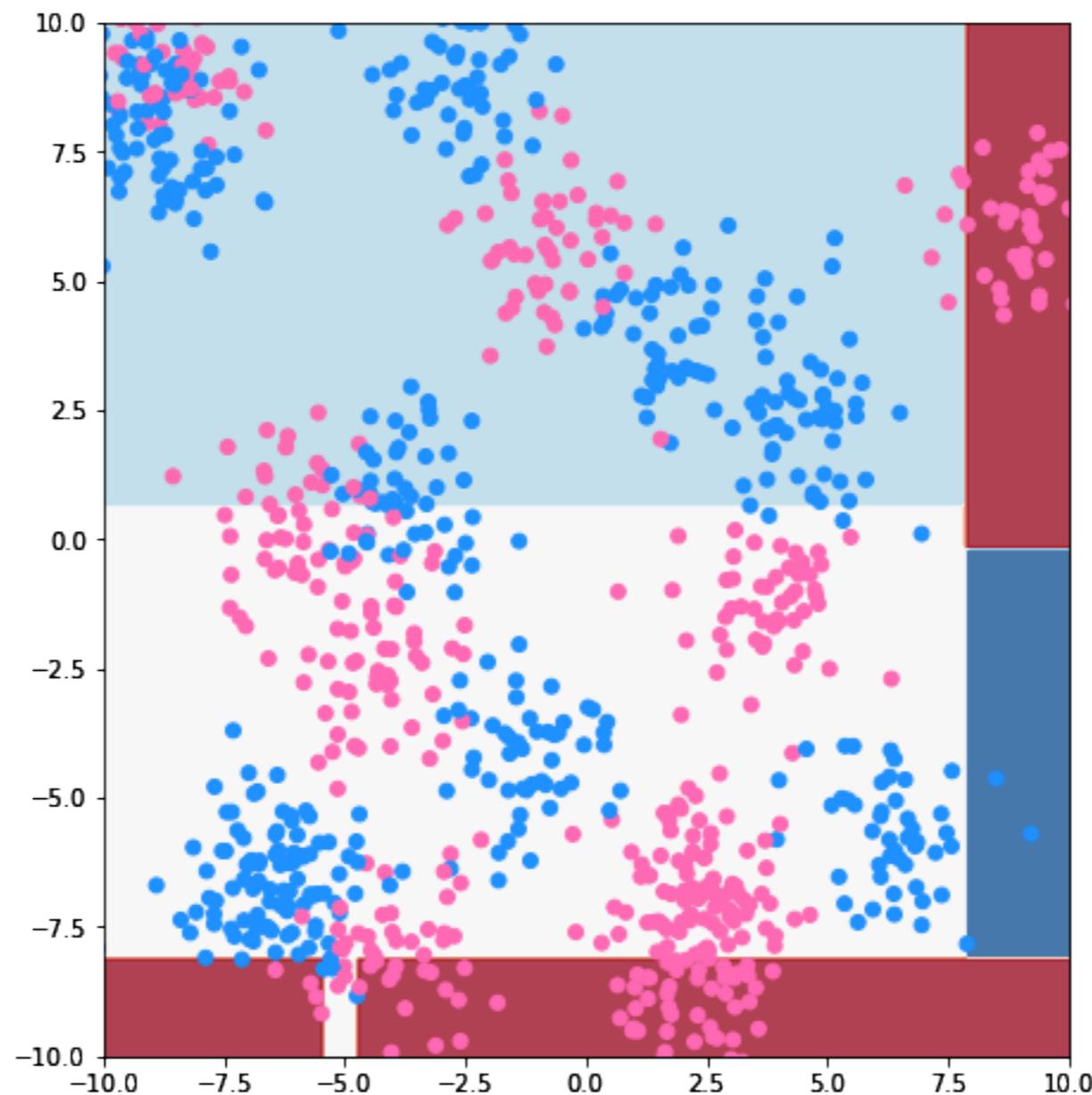
Decision tree, depth=1



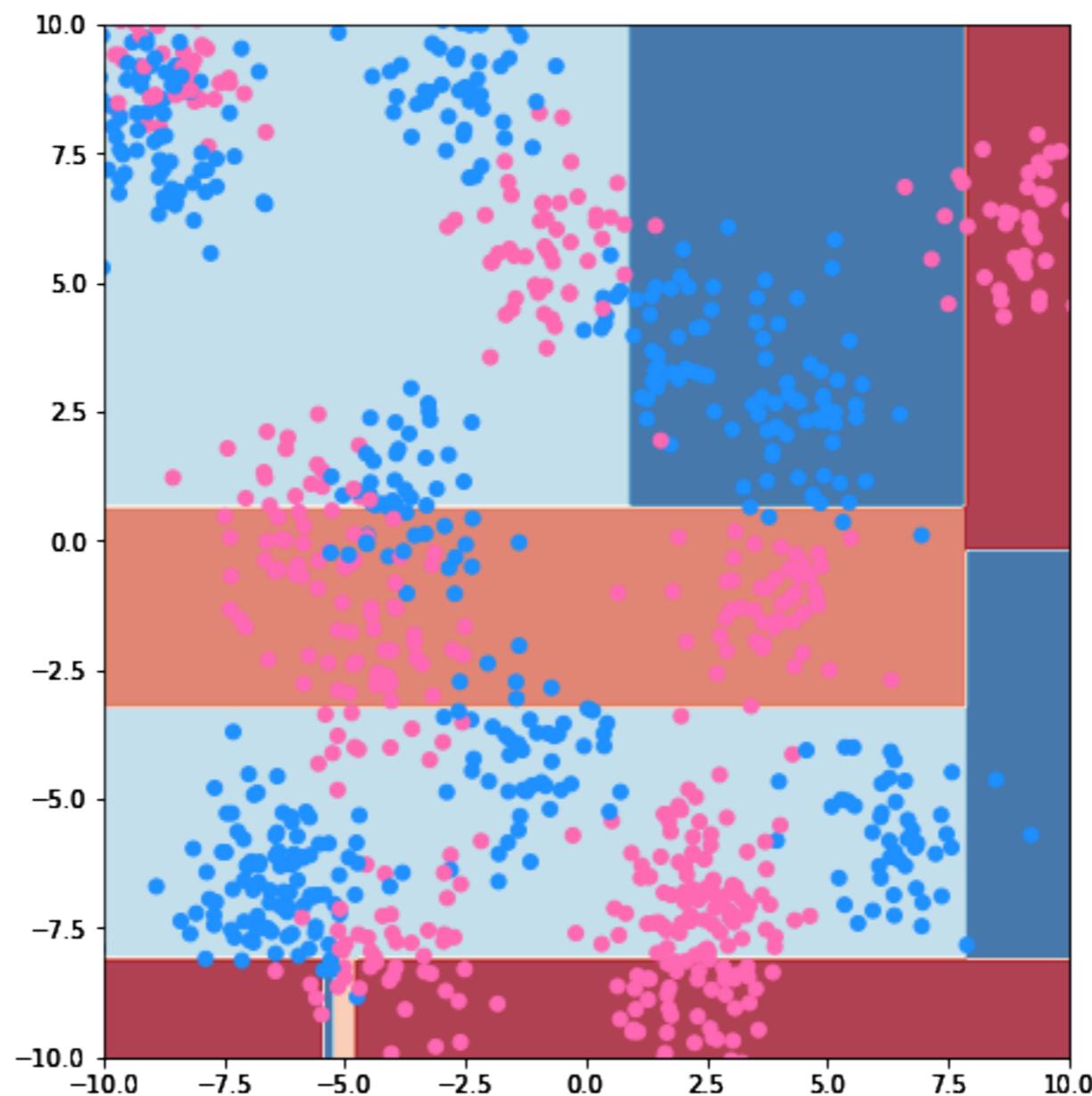
Decision tree, depth=2



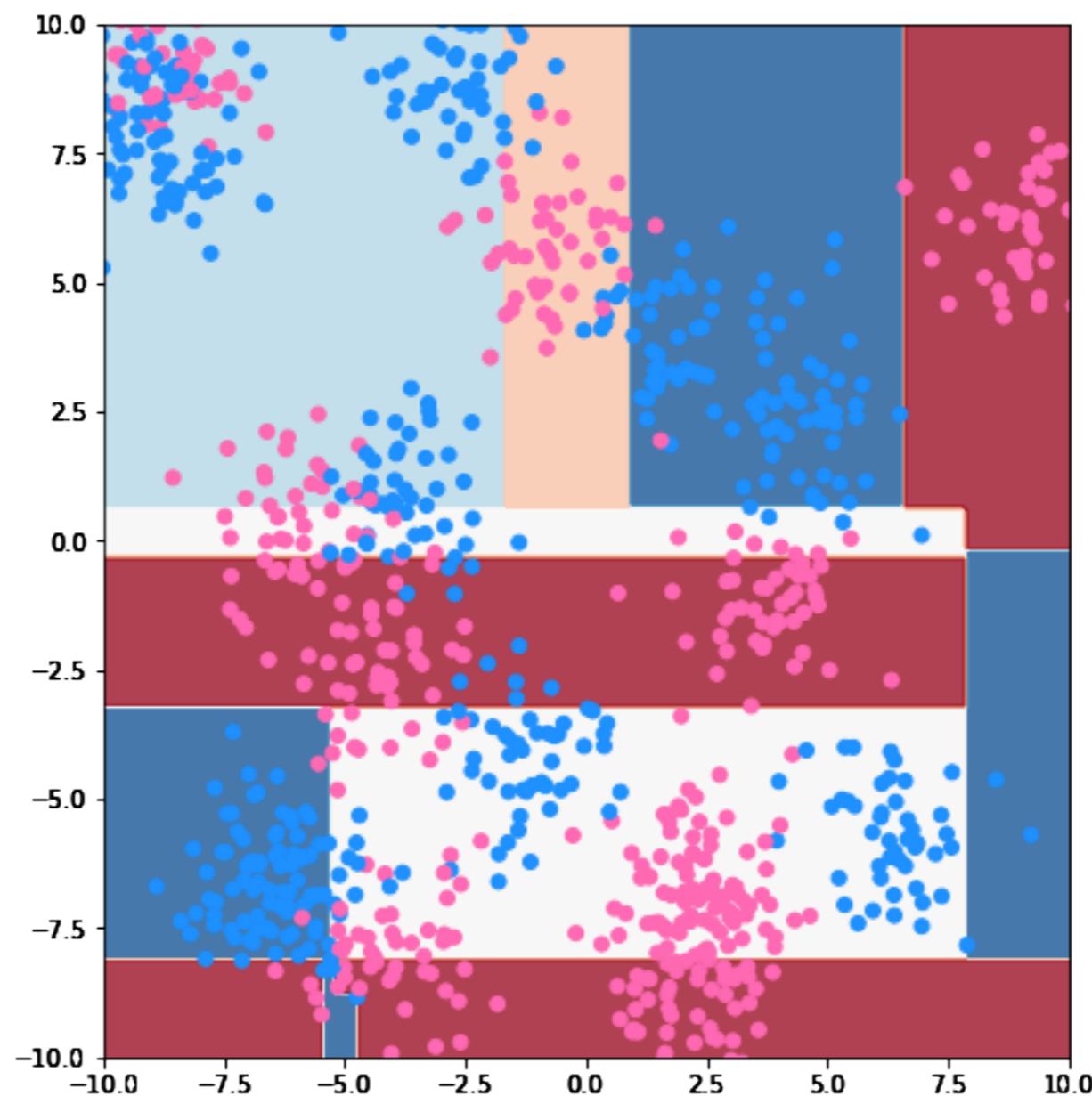
Decision tree, depth=3



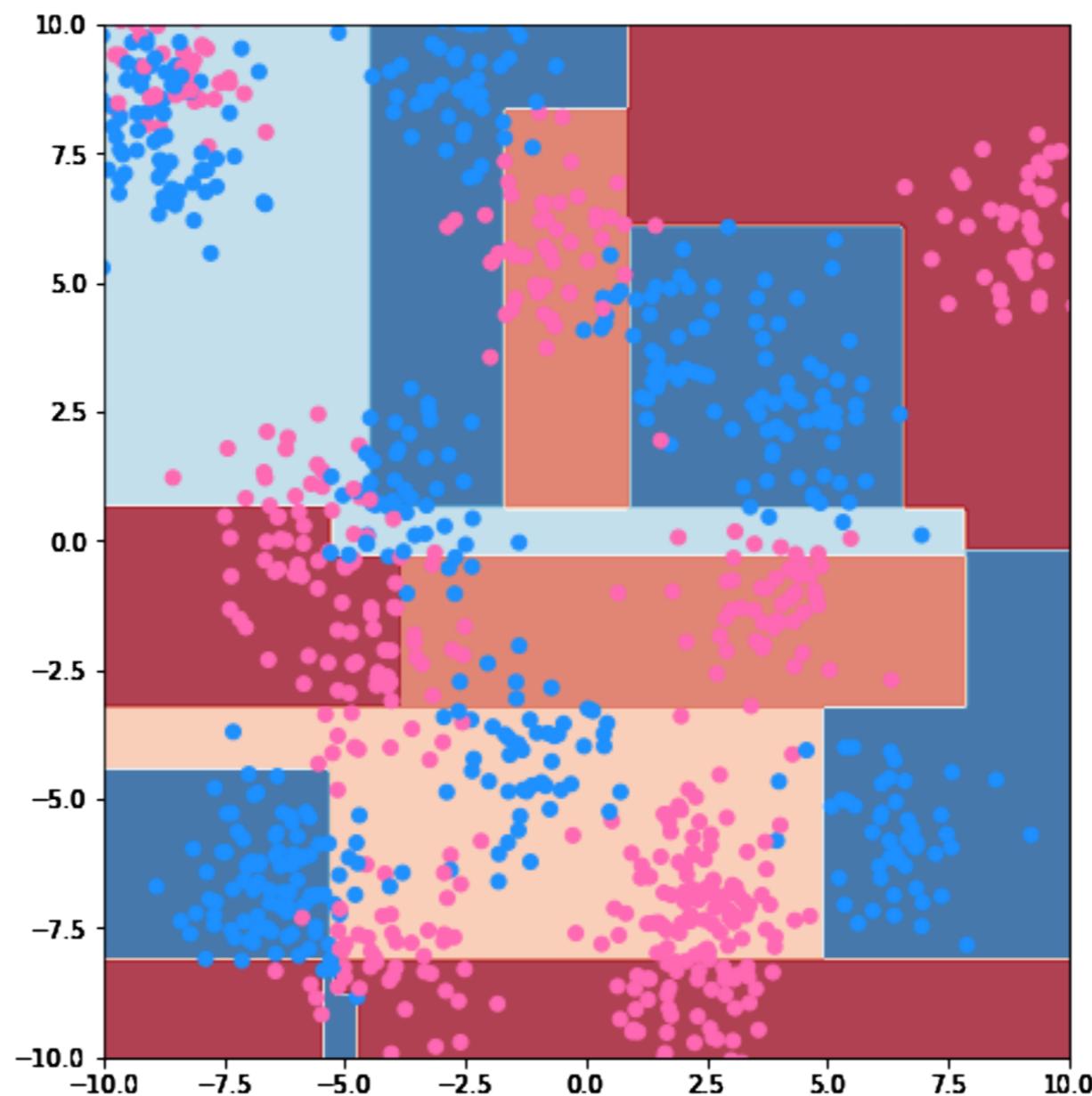
Decision tree, depth=4



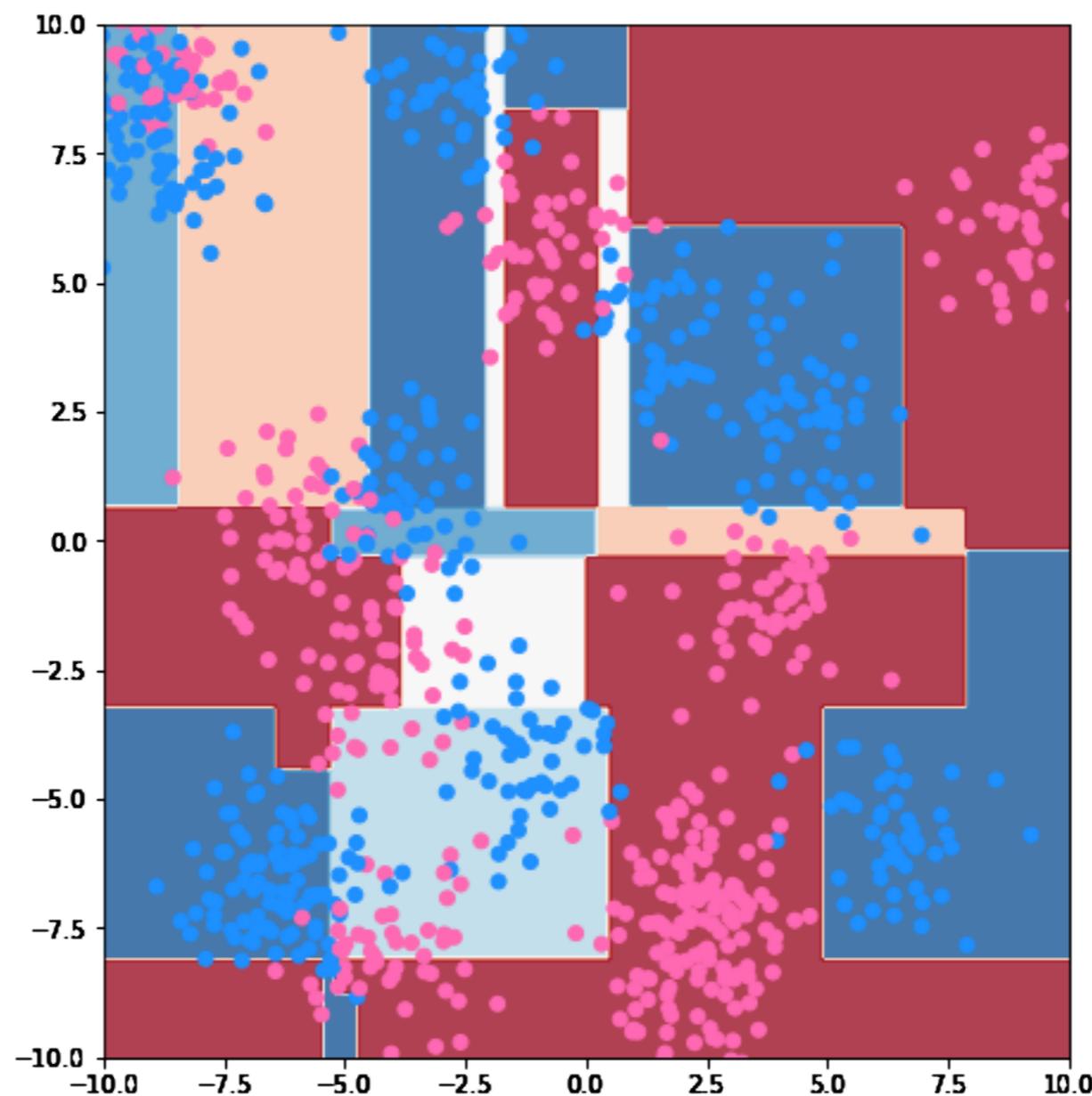
Decision tree, depth=5



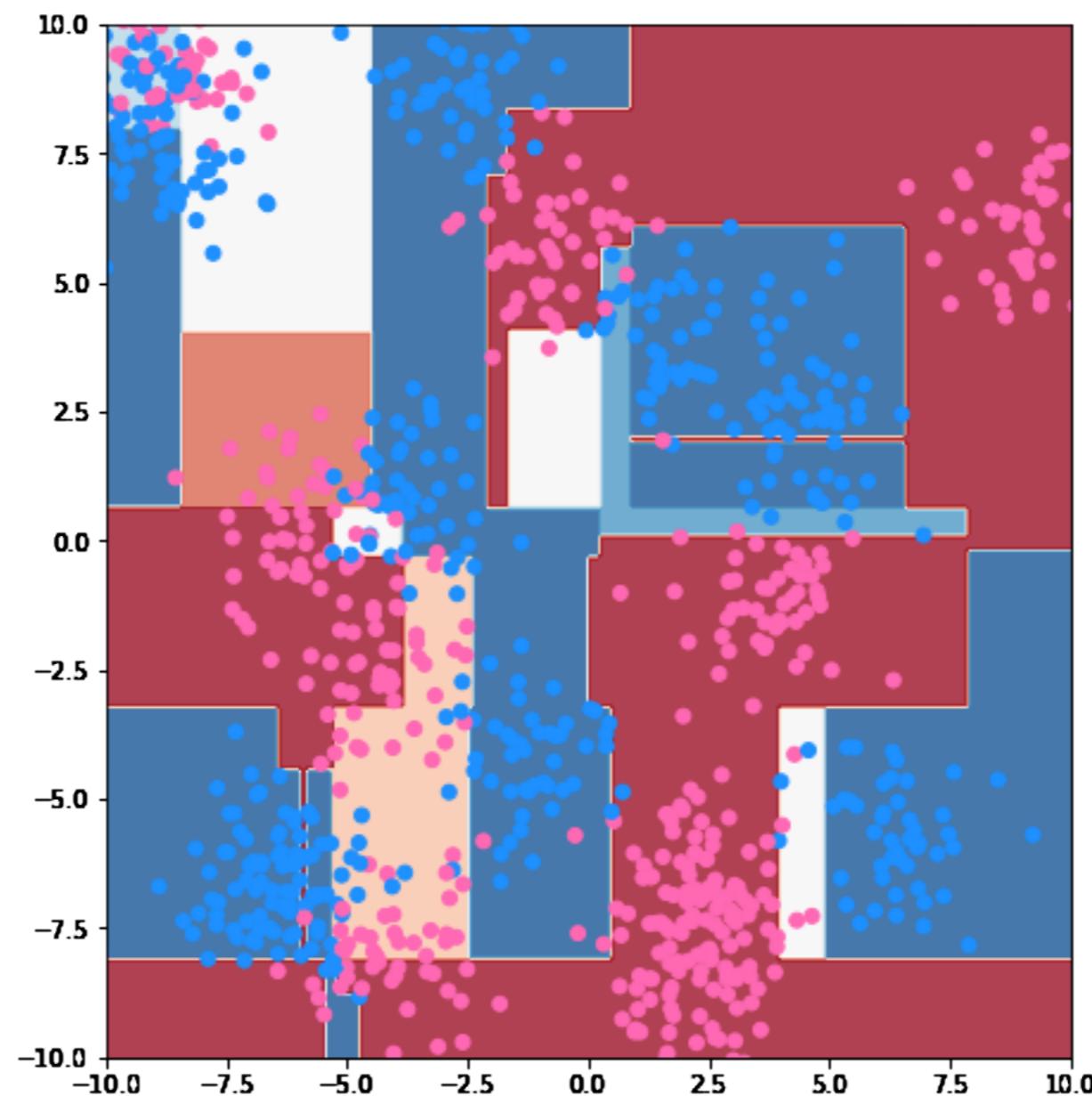
Decision tree, depth=6



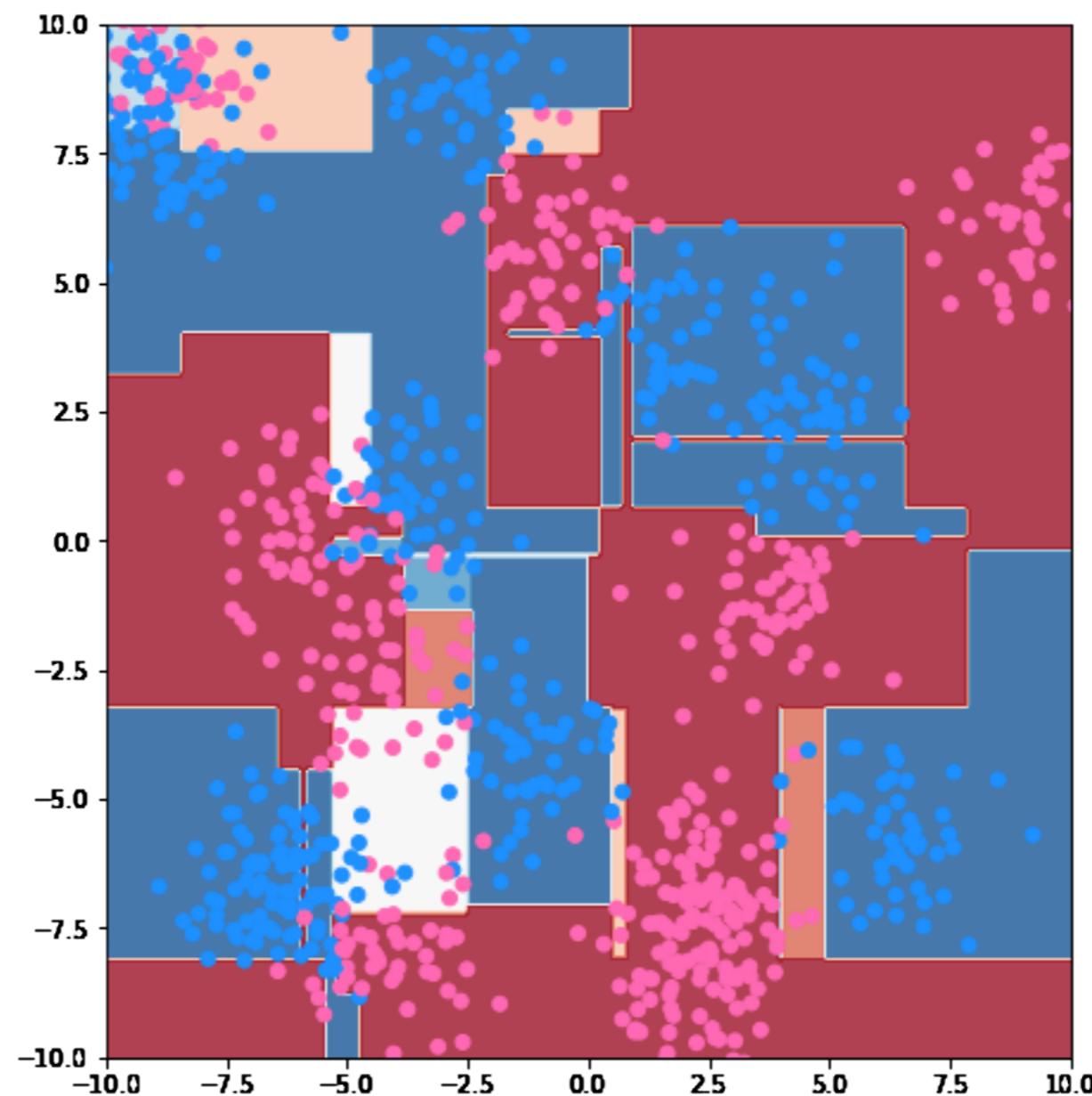
Decision tree, depth=7

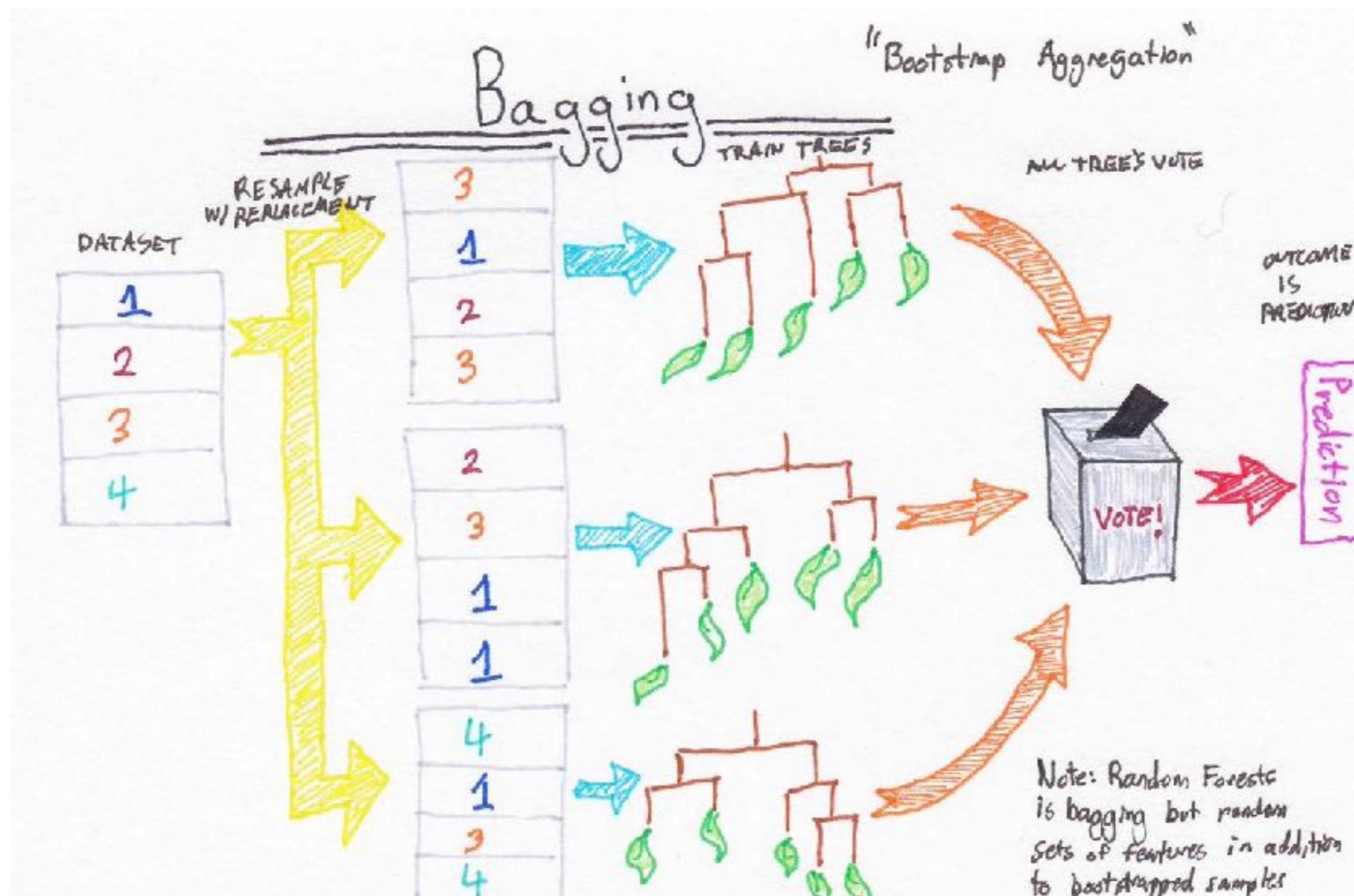


Decision tree, depth=8

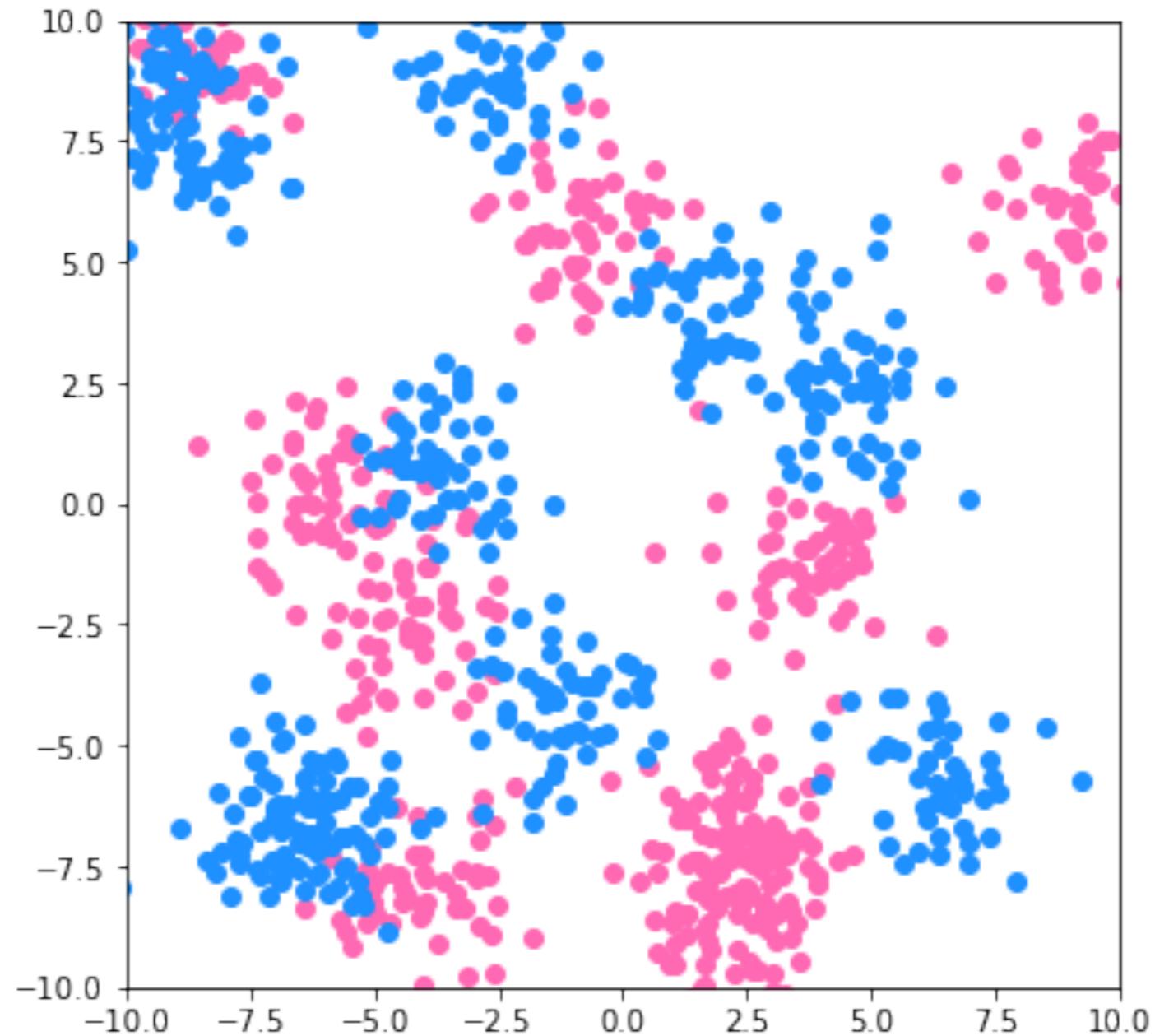


Decision tree, depth=9

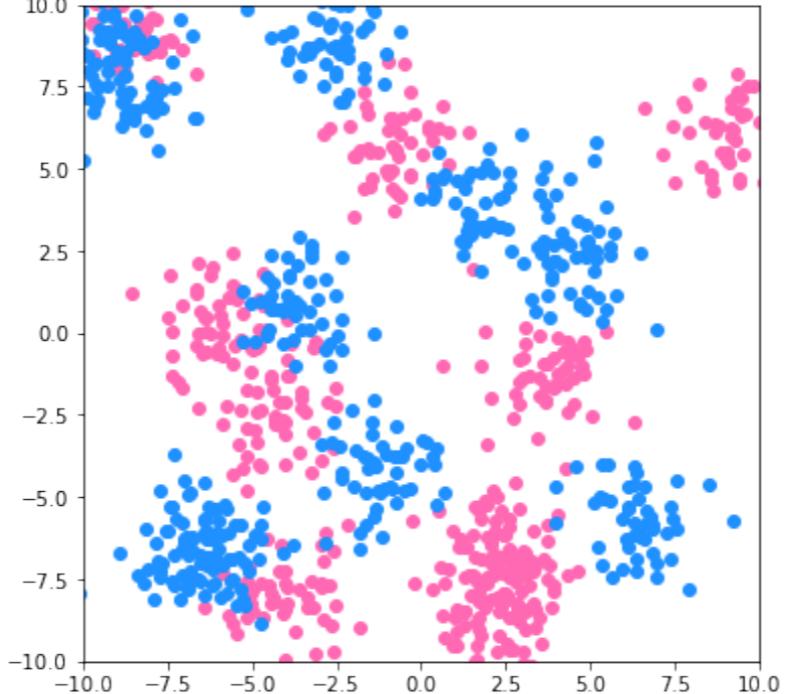




The original set



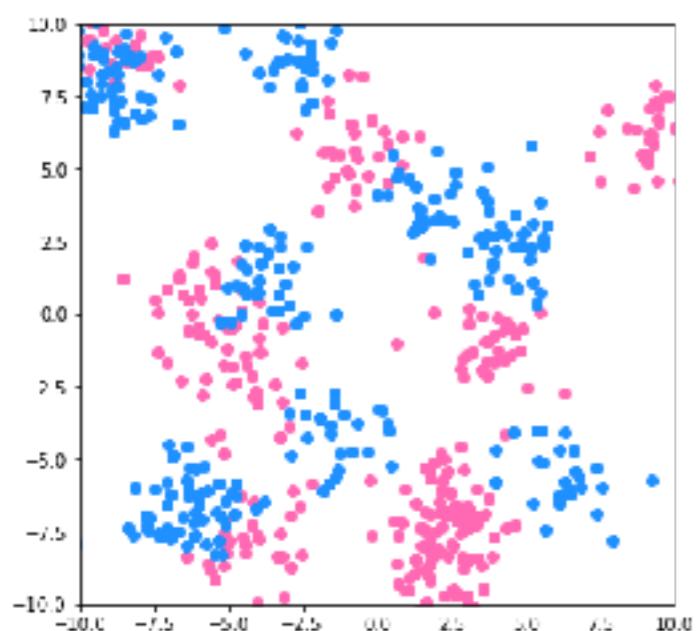
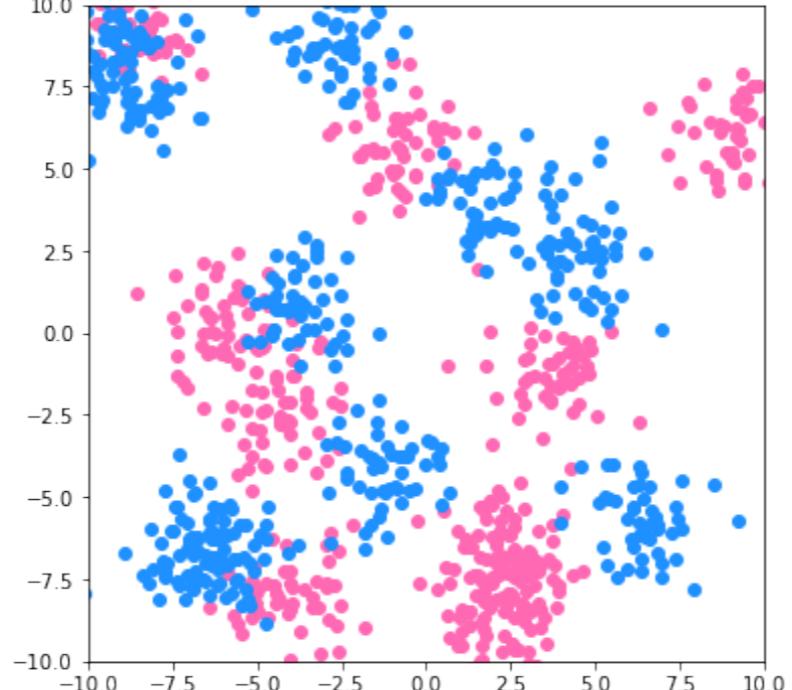
BOOSTRAP!!



BOOSTRAP!!

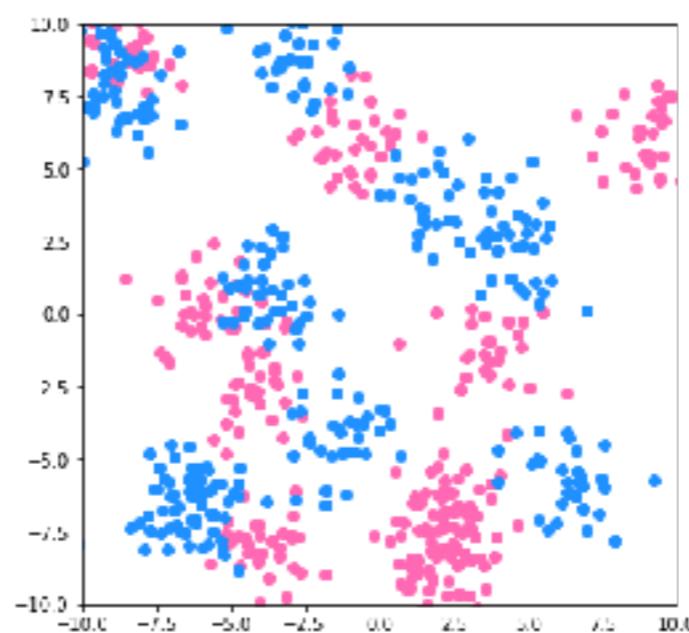
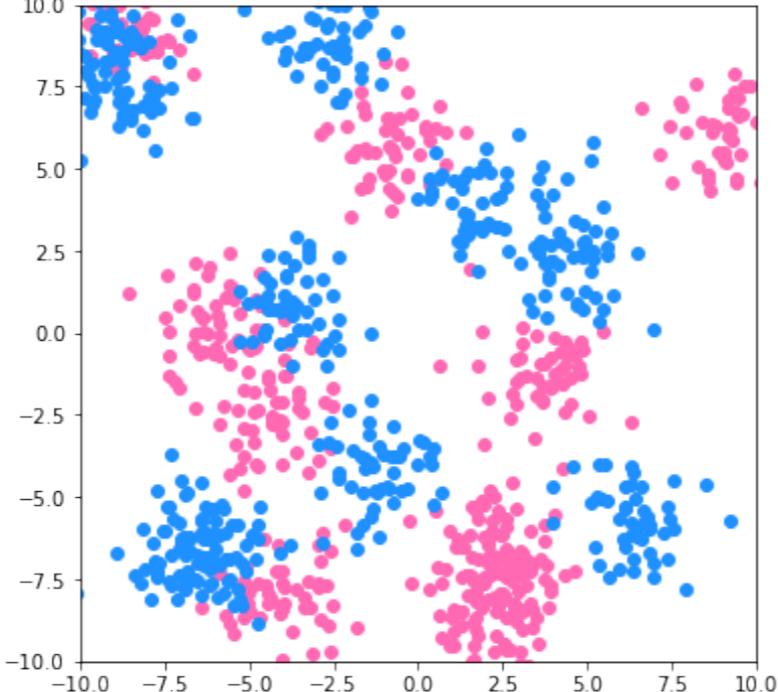
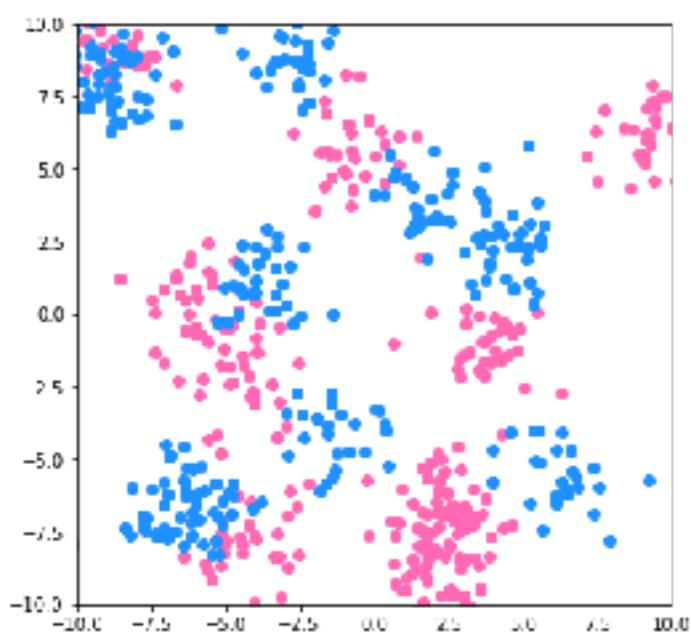
BOOTSTRAP!!

BOOTSTRAP!!



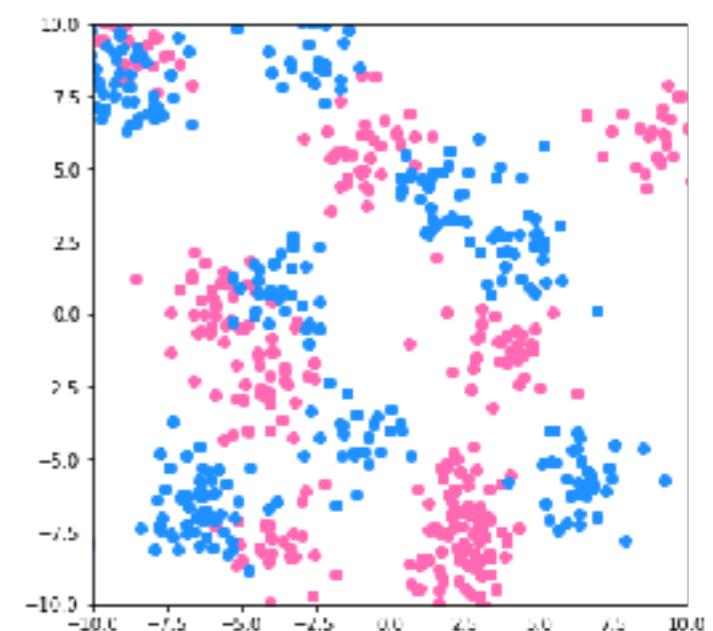
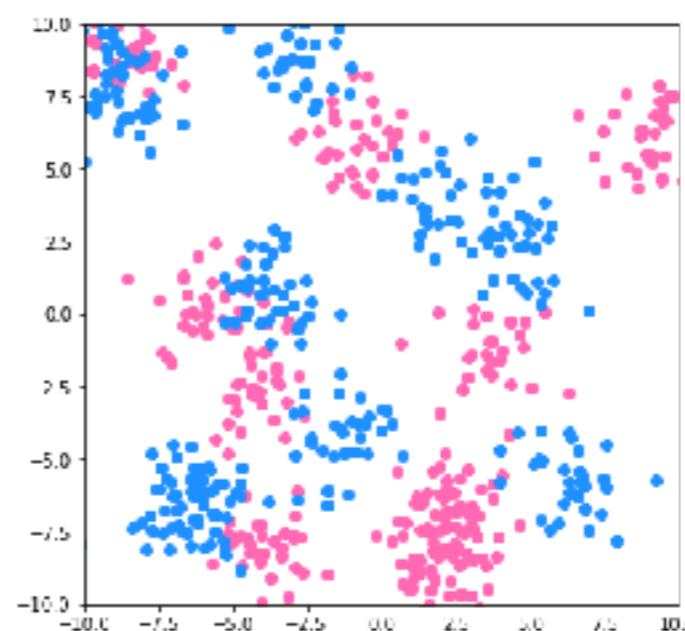
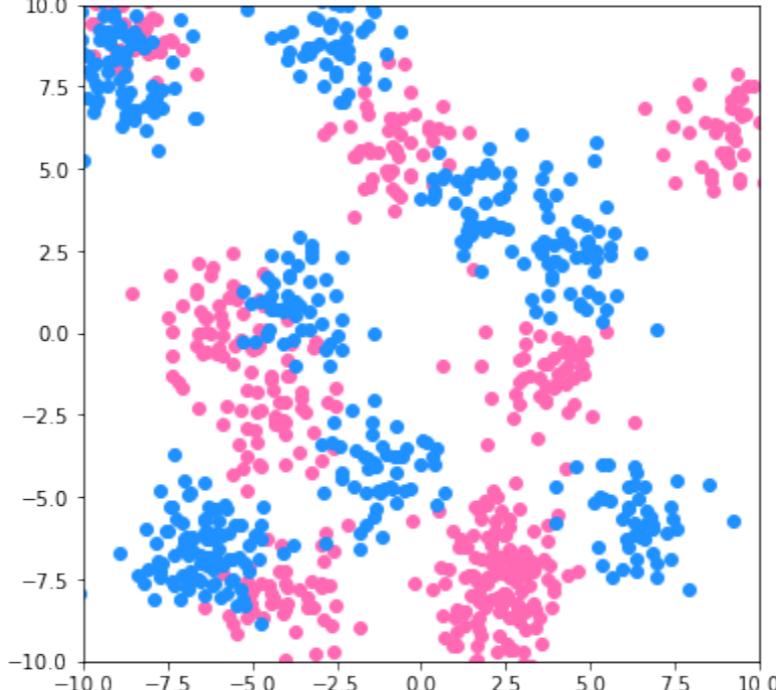
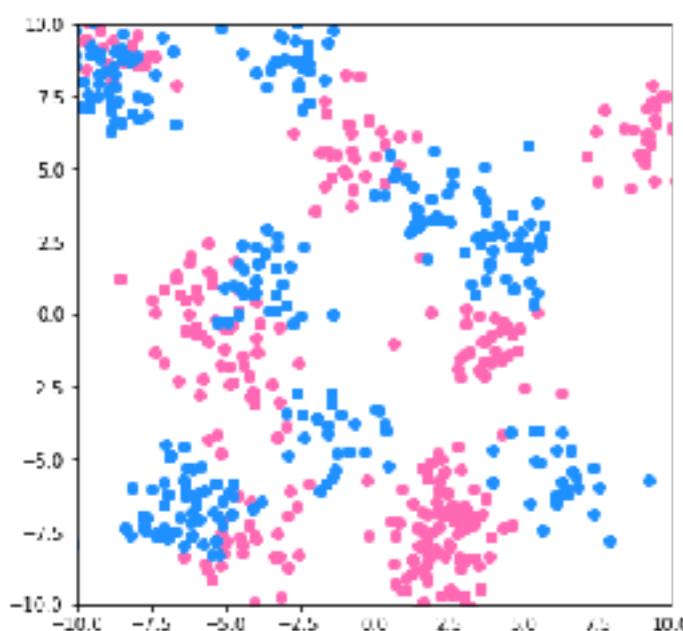
BOOTSTRAP!!

BOOTSTRAP!!



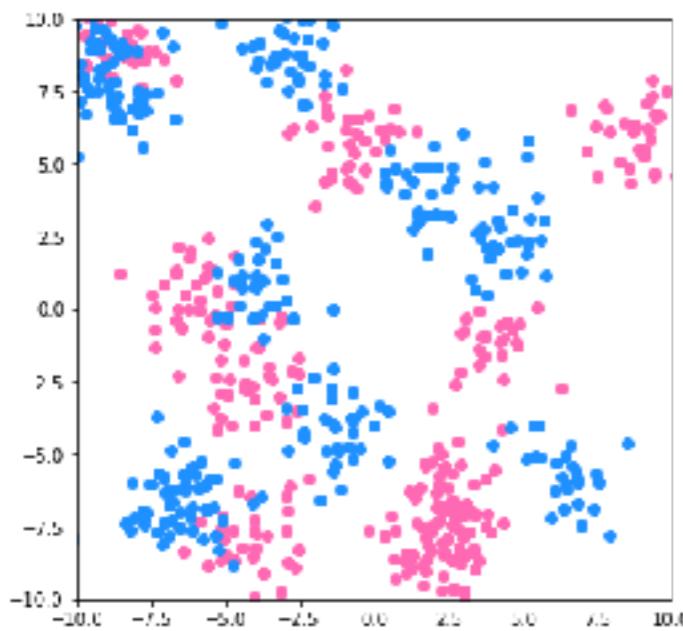
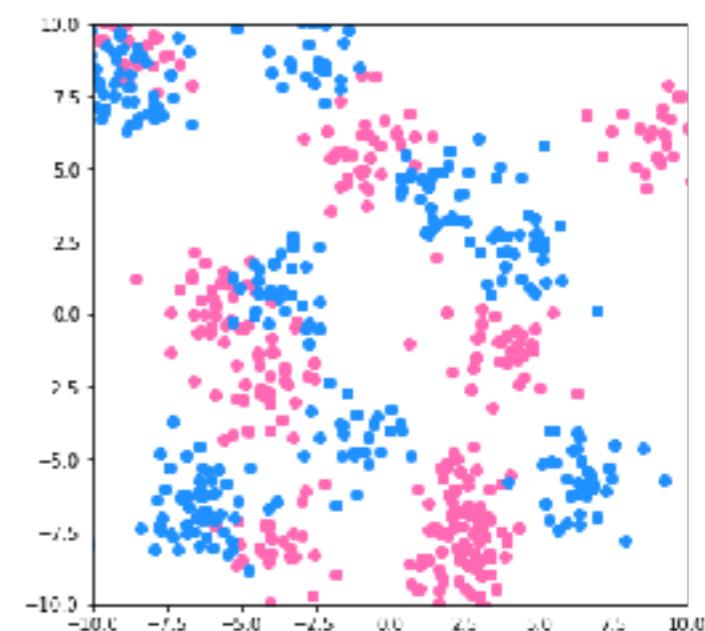
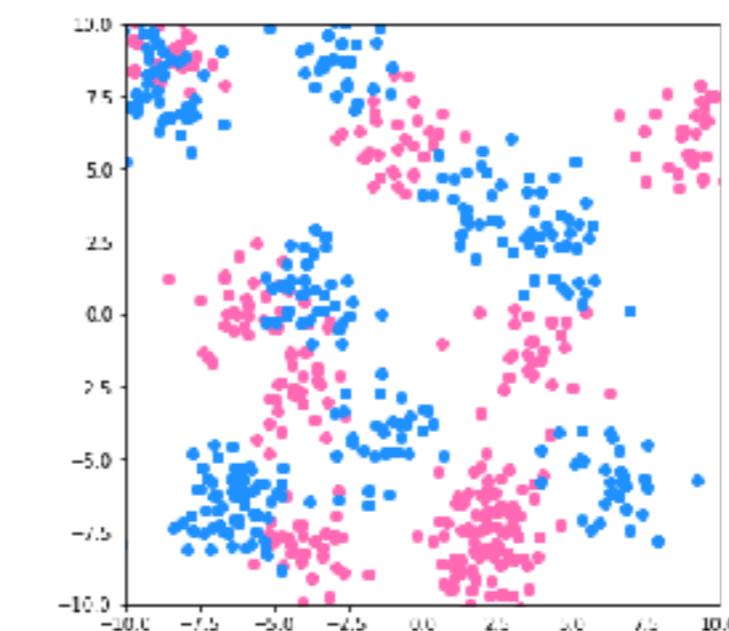
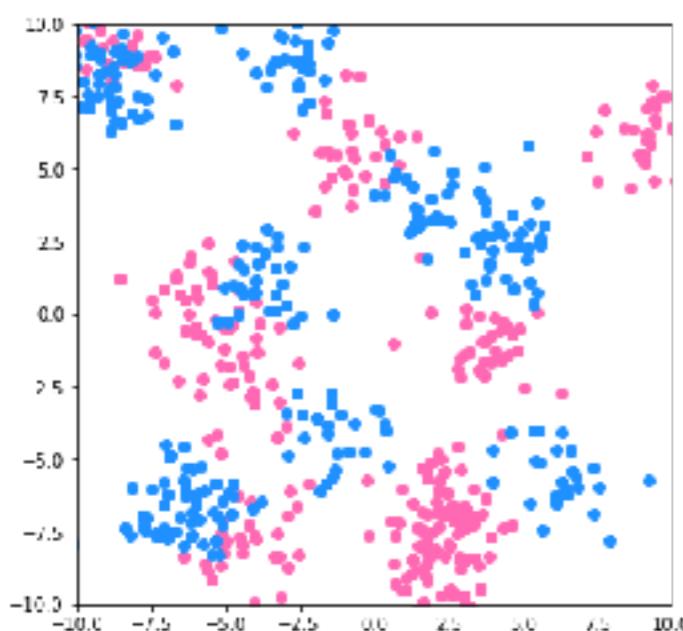
BOOTSTRAP!!

BOOTSTRAP!!



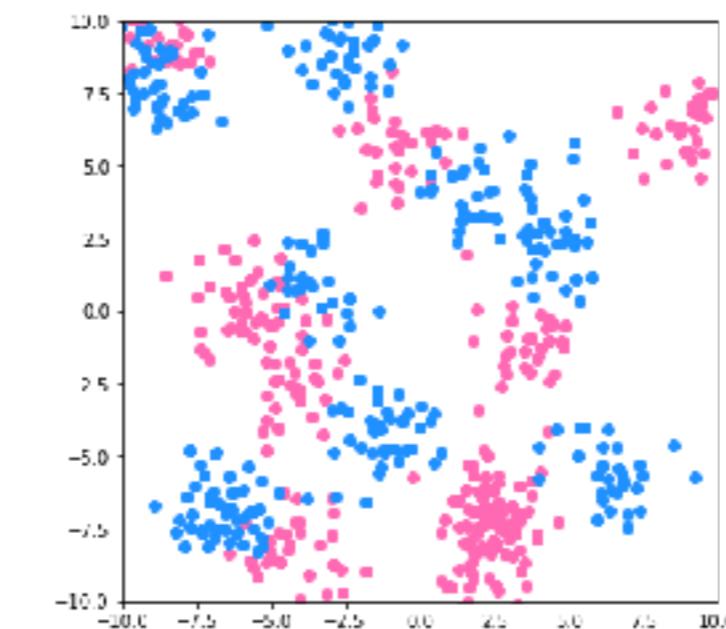
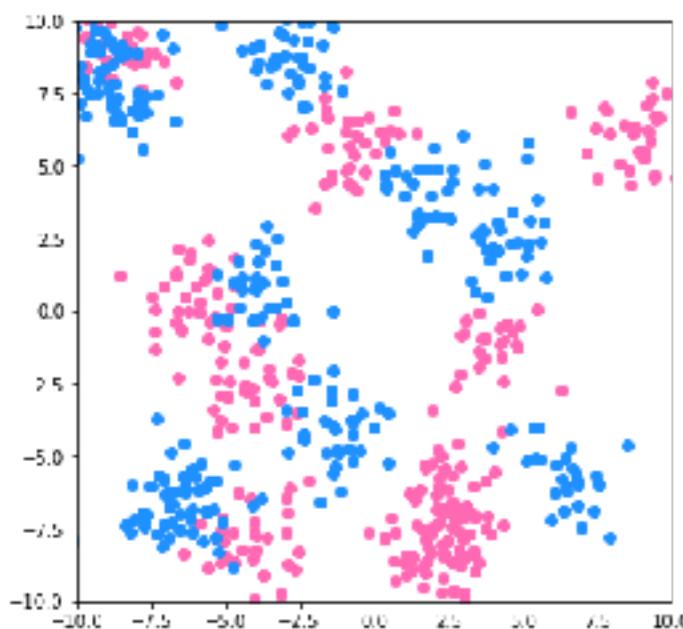
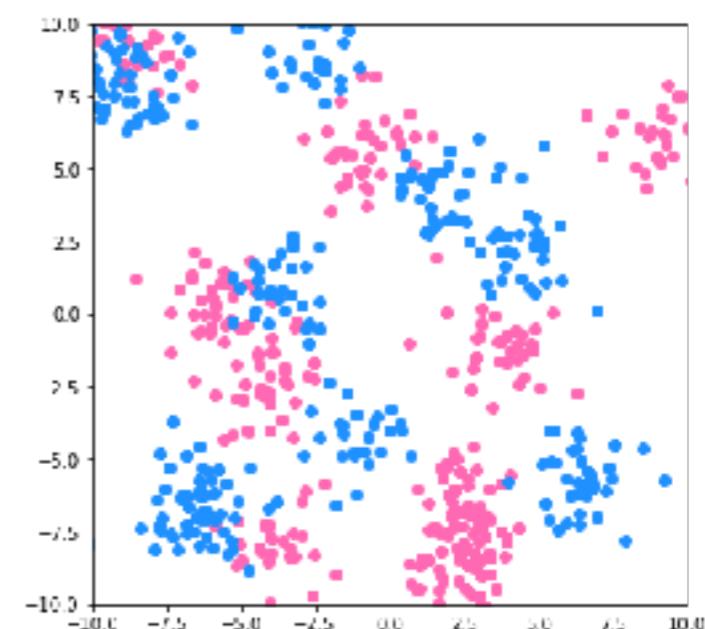
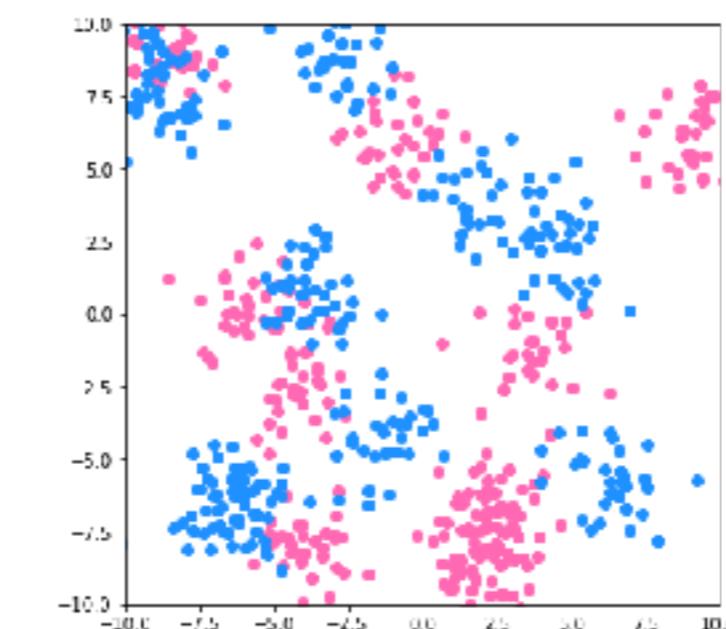
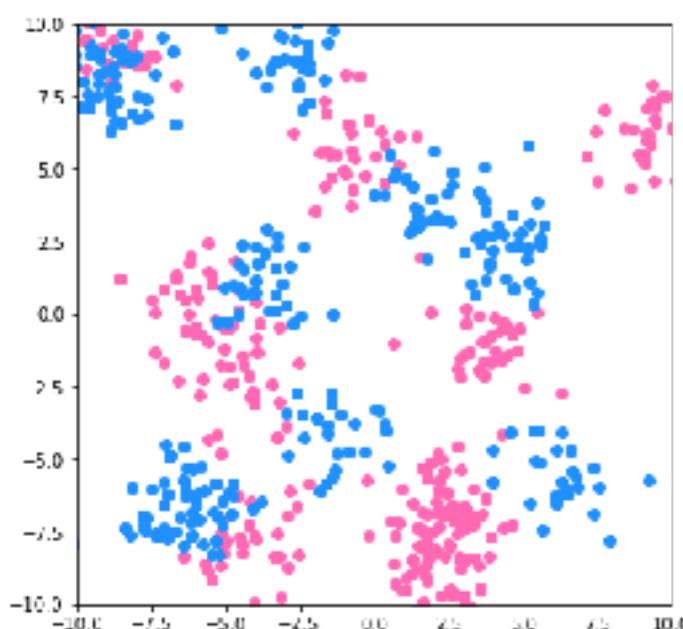
BOOTSTRAP!!

BOOTSTRAP!!



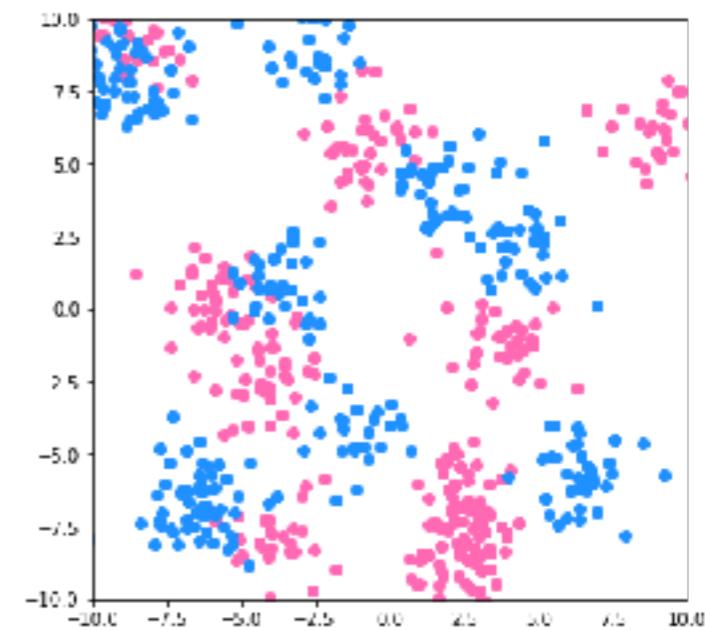
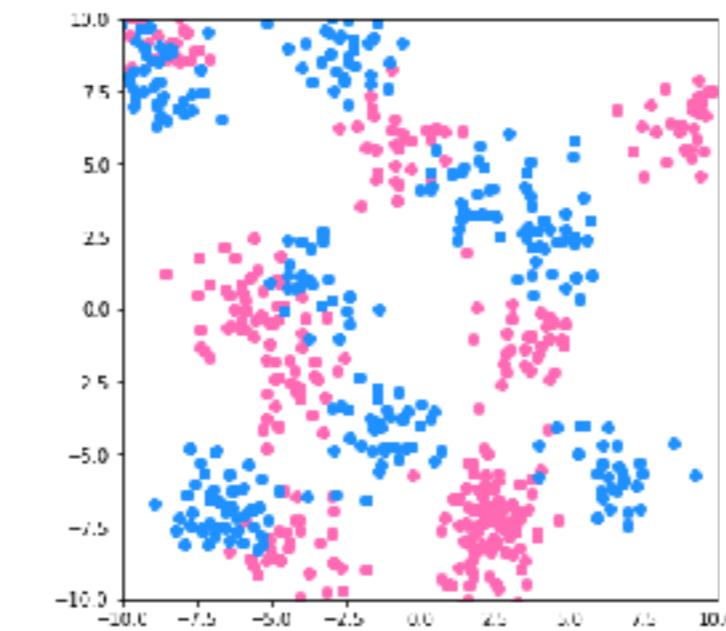
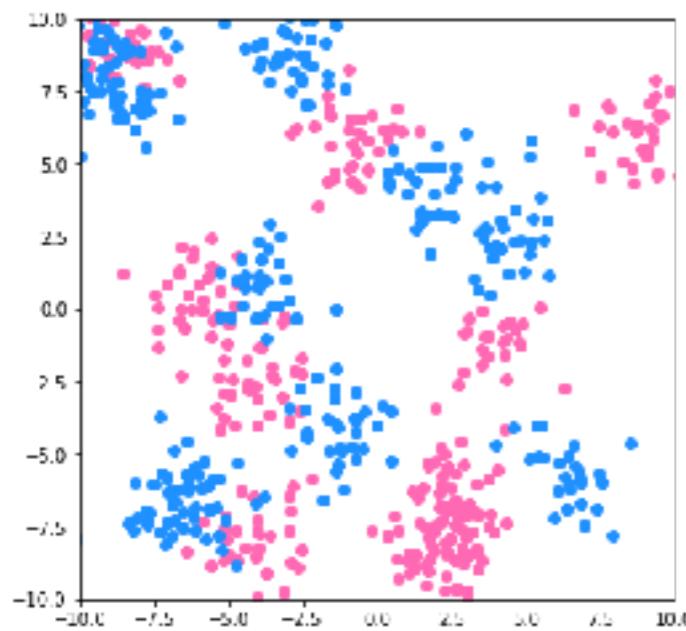
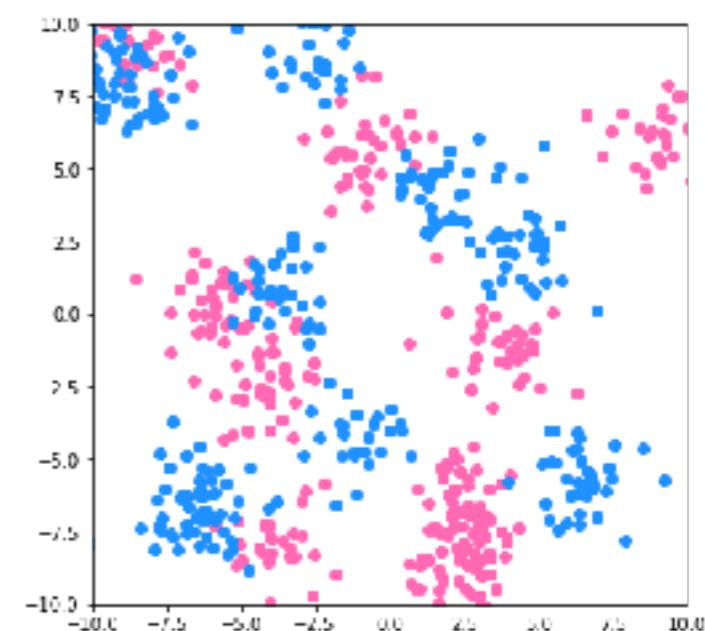
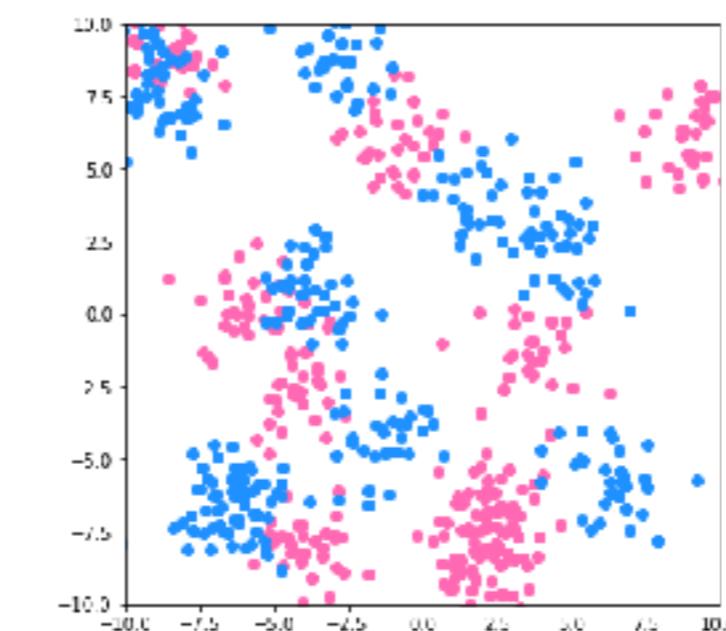
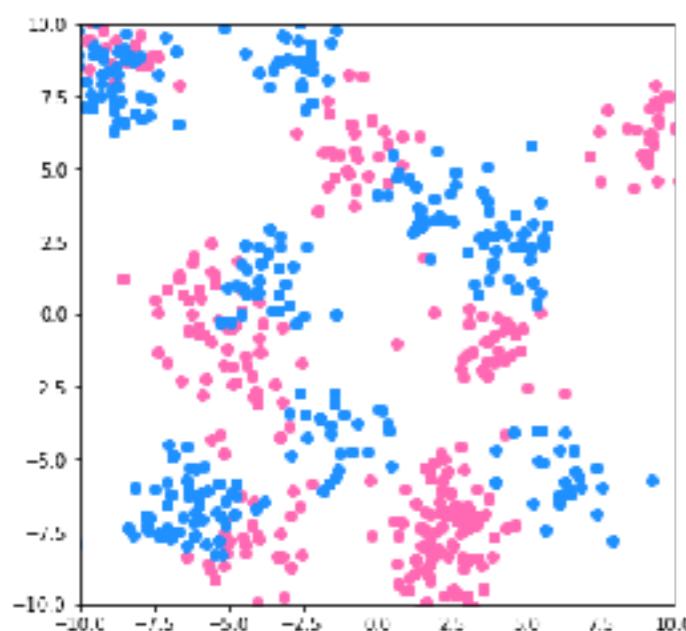
BOOTSTRAP!!

BOOTSTRAP!!

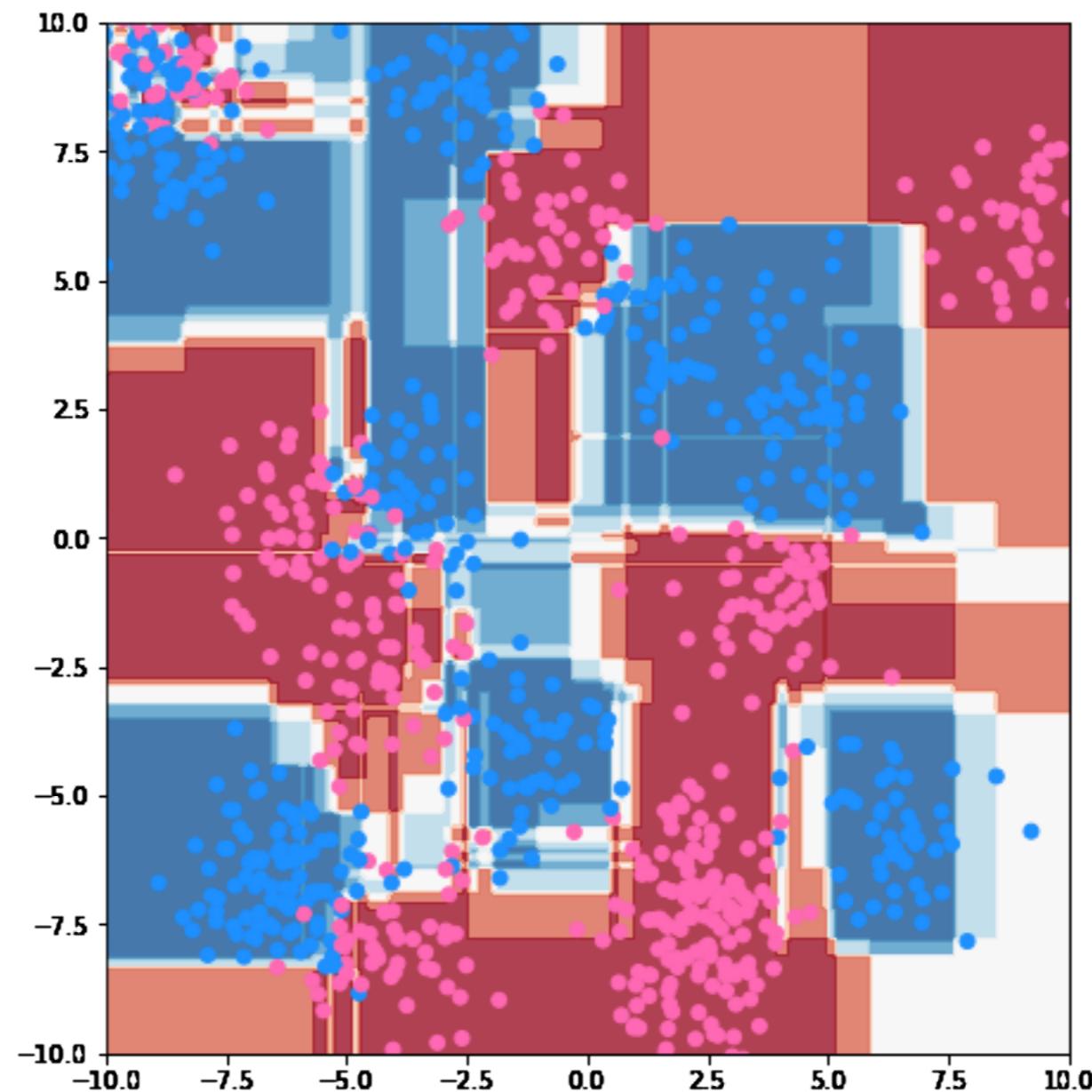


BOOTSTRAP!!

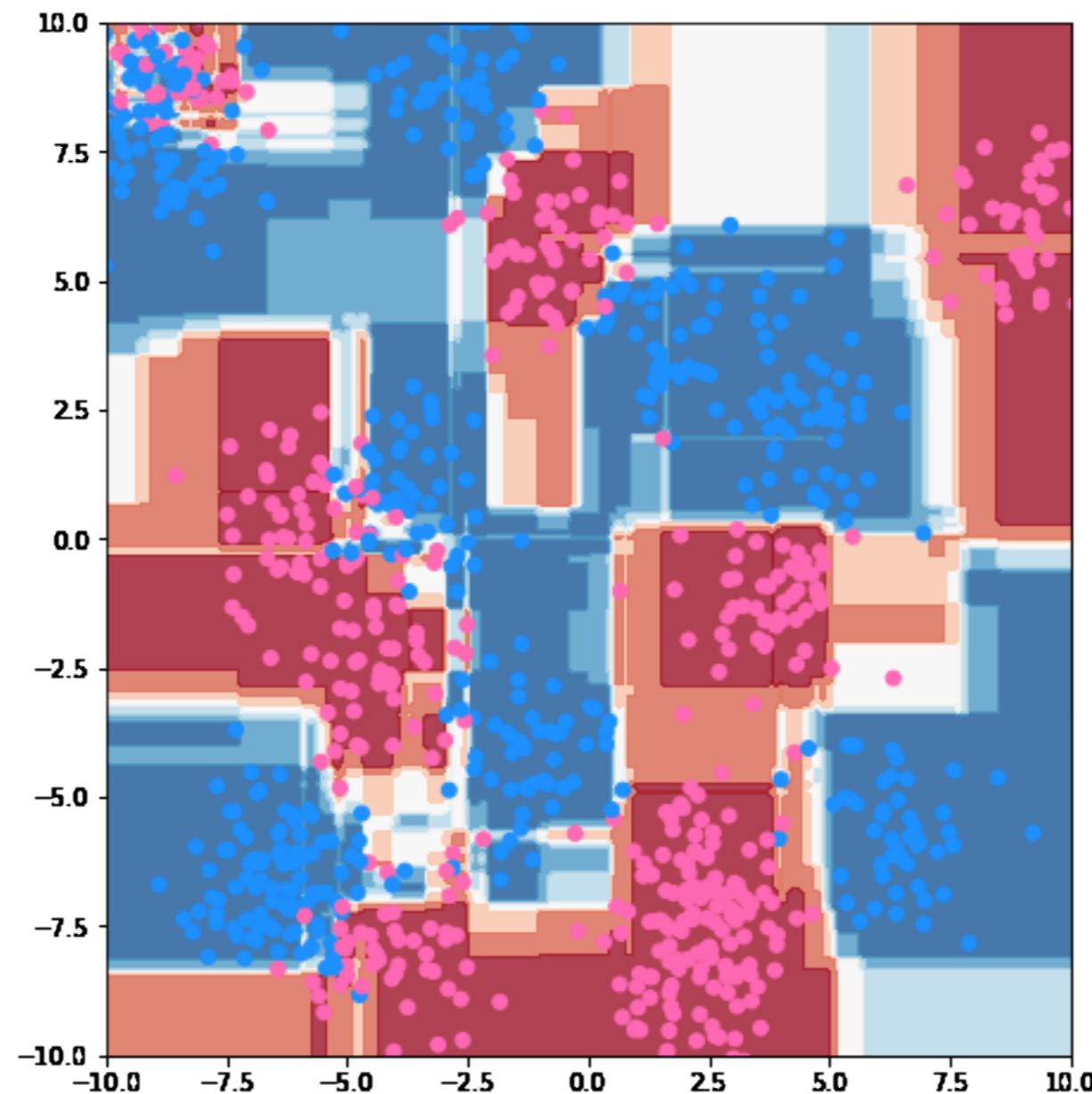
BOOTSTRAP!!



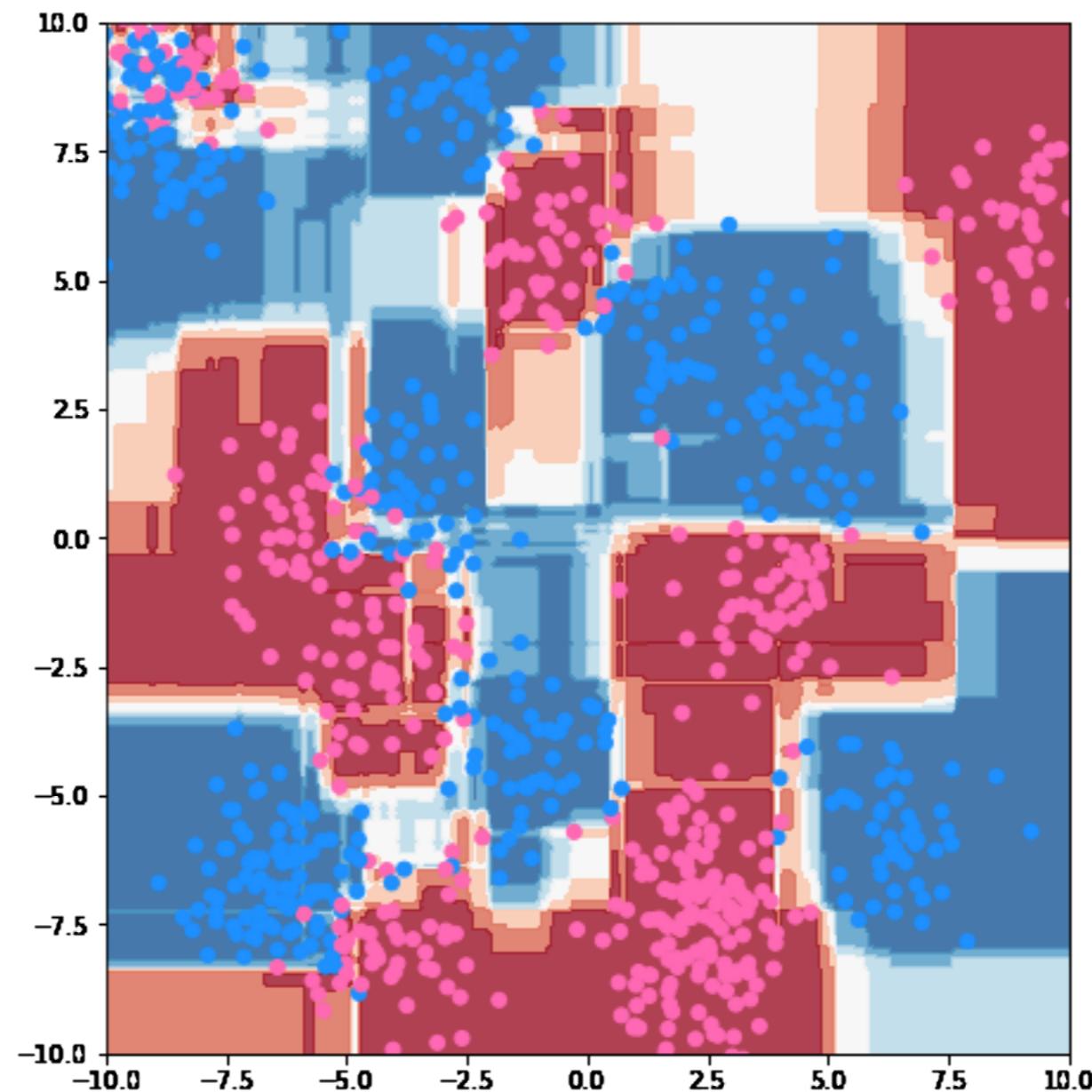
Bagging 5 Decision tree!



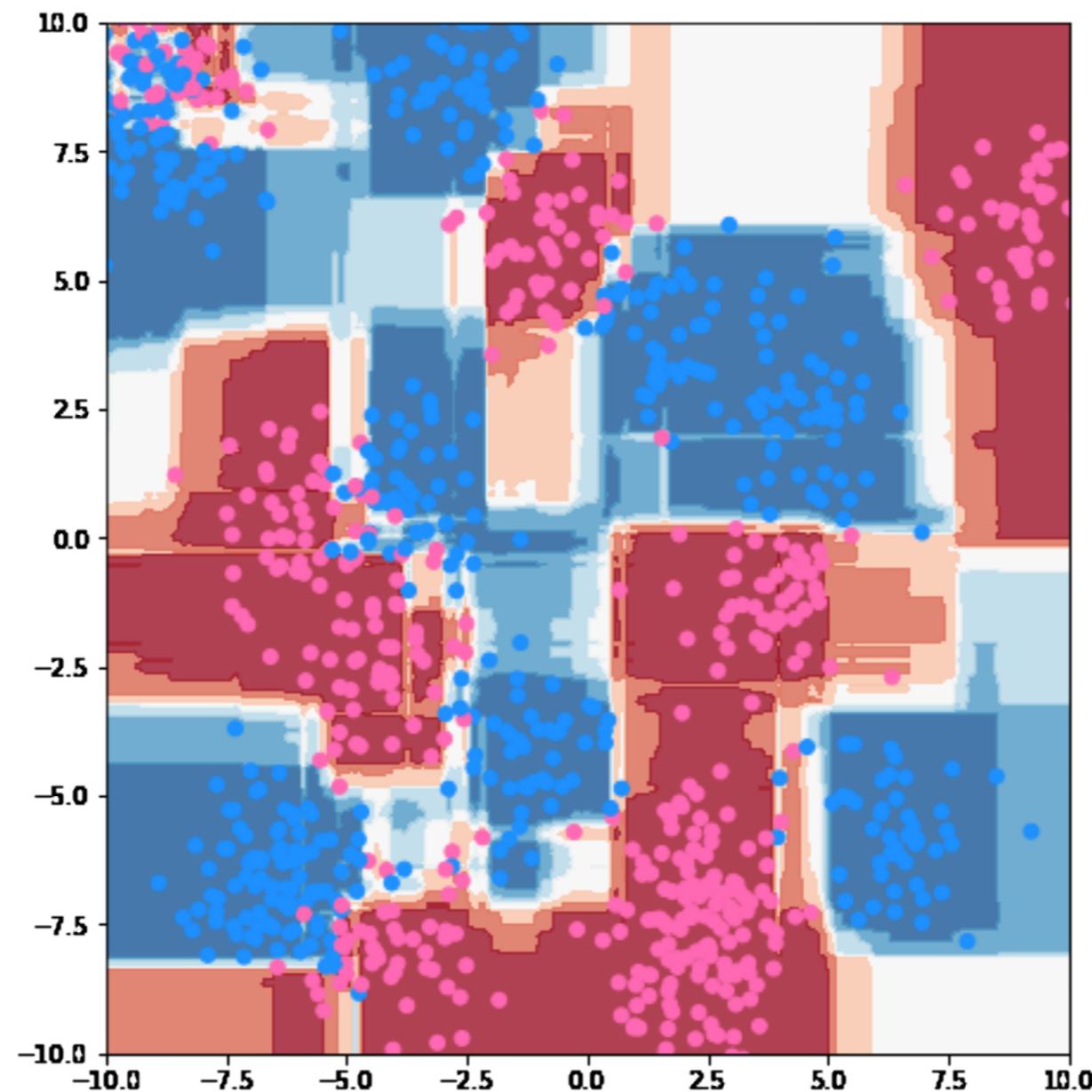
Bagging 10 Decision tree!



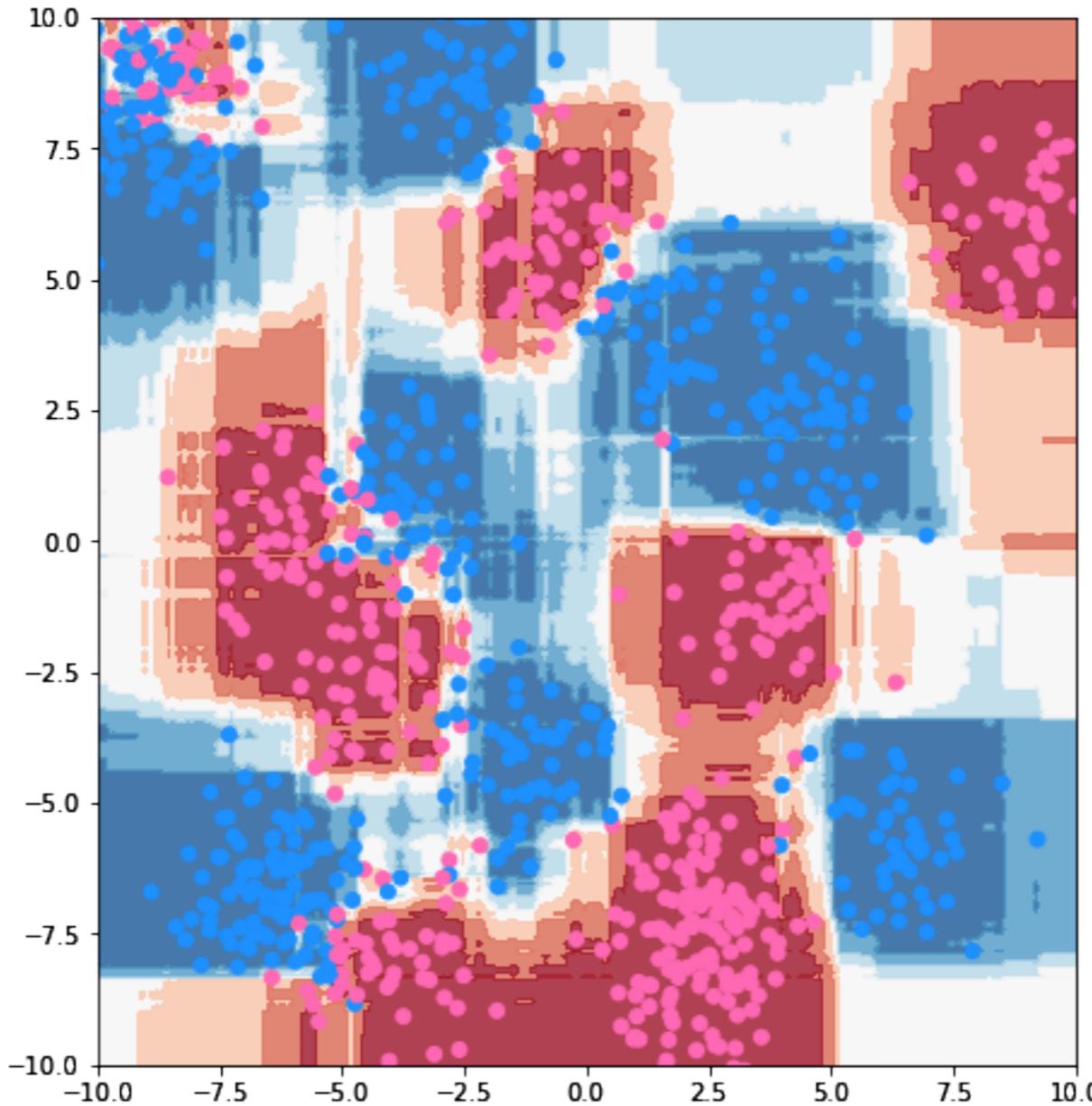
Bagging 20 Decision tree!



Bagging 20 Decision tree!



In practice: Random Forrest



Further randomization: at each candidate split in the learning process, a [random subset of the features](#).
For a classification problem with p features, **\sqrt{p} (rounded down)** features are used in each split

Boosting

Can we make many dumb learners smart?

Boosting

Can we make many dumb learners smart?

Yes (but you need a smart leader!)

Can we make many dumb learners smart?

Can we make many dumb learners smart?

Kearns and Valiant '88:

- Does weak learnability imply strong learnability?
In other words, can we boost from weak to strong?

Can we make many dumb learners smart?

Kearns and Valiant '88:

- Does weak learnability imply strong learnability?
In other words, can we boost from weak to strong?

Schapire '89

Freund and Schapire '95

- Yes, with adaBoost

How to combine many weak classifiers?

Ingredients:

- (i) T classifiers $h_t(\cdot)$, each of them slightly better than random
- (ii) A training set with N labeled examples

Adaptive linear combination of classifiers

ADABOOST

$$H(\vec{x}) = \sum_{t=1}^T \alpha_t h_t(\vec{x})$$

Exponential loss:

$$\mathcal{R} = \sum_i e^{-Y_i H(\vec{x}_i)}$$

Goal: start from $t=0$, add new classifiers and

- (i) Adapt the weight alpha at each steps
- (ii) re-weight the instances in the training set : more weight to ill-classified instances
(each new classifier concentrate on badly classified examples)

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i e^{-Y_i H(\vec{x}_i)}$$

$$H(\vec{x}) = \sum_{t=1}^T \alpha_t h_t(\vec{x})$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i e^{-Y_i H(\vec{x}_i)} \quad H(\vec{x}) = \sum_{t=1}^T \alpha_t h_t(\vec{x})$$

$$\mathcal{R} = \sum_i e^{-Y_i \sum_{t=1}^{\tau-1} \alpha_t h_t(\vec{x}_i) - Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i e^{-Y_i H(\vec{x}_i)}$$
$$H(\vec{x}) = \sum_{t=1}^T \alpha_t h_t(\vec{x})$$
$$\mathcal{R} = \sum_i e^{-Y_i \underbrace{\sum_{t=1}^{\tau-1} \alpha_t h_t(\vec{x}_i)}_{\lambda_i^\tau} - Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i e^{-Y_i H(\vec{x}_i)} \quad H(\vec{x}) = \sum_{t=1}^T \alpha_t h_t(\vec{x})$$

$$\mathcal{R} = \sum_i e^{-Y_i \underbrace{\sum_{t=1}^{\tau-1} \alpha_t h_t(\vec{x}_i)}_{\lambda_i^\tau} - Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = \sum_i e^{-Y_i \lambda_i^\tau} e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i e^{-Y_i H(\vec{x}_i)} \quad H(\vec{x}) = \sum_{t=1}^T \alpha_t h_t(\vec{x})$$

$$\mathcal{R} = \sum_i e^{-Y_i \underbrace{\sum_{t=1}^{\tau-1} \alpha_t h_t(\vec{x}_i)}_{\lambda_i^\tau} - Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = \sum_i e^{-Y_i \underbrace{\lambda_i^\tau}_{\omega_i^\tau}} e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i e^{-Y_i H(\vec{x}_i)} \quad H(\vec{x}) = \sum_{t=1}^T \alpha_t h_t(\vec{x})$$

$$\mathcal{R} = \sum_i e^{-Y_i \underbrace{\sum_{t=1}^{\tau-1} \alpha_t h_t(\vec{x}_i)}_{\lambda_i^\tau} - Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = \sum_i e^{-Y_i \underbrace{\lambda_i^\tau}_{\omega_i^\tau}} e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i e^{-Y_i H(\vec{x}_i)} \quad H(\vec{x}) = \sum_{t=1}^T \alpha_t h_t(\vec{x})$$

$$\mathcal{R} = \sum_i e^{-Y_i \underbrace{\sum_{t=1}^{\tau-1} \alpha_t h_t(\vec{x}_i)}_{\lambda_i^\tau} - Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = \sum_i e^{-Y_i \underbrace{\lambda_i^\tau}_{\omega_i^\tau}} e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

ω_i^τ = weights for each instances at time τ

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_{i \in \text{OK}} \omega_i^\tau + e^{\alpha_\tau} \sum_{i \in \text{NOT OK}} \omega_i^\tau$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_{i \in \text{OK}} \omega_i^\tau + e^{\alpha_\tau} \sum_{i \in \text{NOT OK}} \omega_i^\tau = e^{-\alpha_\tau} \sum_i \mathbf{1}(Y_i = h_i) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_{i \in \text{OK}} \omega_i^\tau + e^{\alpha_\tau} \sum_{i \in \text{NOT OK}} \omega_i^\tau = e^{-\alpha_\tau} \sum_i \mathbf{1}(Y_i = h_i) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_i (1 - \mathbf{1}(Y_i \neq h_i)) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_{i \in \text{OK}} \omega_i^\tau + e^{\alpha_\tau} \sum_{i \in \text{NOT OK}} \omega_i^\tau = e^{-\alpha_\tau} \sum_i \mathbf{1}(Y_i = h_i) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_i (1 - \underbrace{\mathbf{1}(Y_i \neq h_i)}_{\epsilon_t}) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_{i \in \text{OK}} \omega_i^\tau + e^{\alpha_\tau} \sum_{i \in \text{NOT OK}} \omega_i^\tau = e^{-\alpha_\tau} \sum_i \mathbf{1}(Y_i = h_i) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_i (1 - \underbrace{\mathbf{1}(Y_i \neq h_i)}_{\epsilon_t}) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

ϵ_t is what the classifier $h_t(\cdot)$ is trying to minimise

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_{i \in \text{OK}} \omega_i^\tau + e^{\alpha_\tau} \sum_{i \in \text{NOT OK}} \omega_i^\tau = e^{-\alpha_\tau} \sum_i \mathbf{1}(Y_i = h_i) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_i (1 - \underbrace{\mathbf{1}(Y_i \neq h_i)}_{\epsilon_t}) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

ϵ_t is what the classifier $h_t(\cdot)$ is trying to minimise

Let us choose α in order to minimize the global risk

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau h_\tau(\vec{x}_i)}$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_{i \in \text{OK}} \omega_i^\tau + e^{\alpha_\tau} \sum_{i \in \text{NOT OK}} \omega_i^\tau = e^{-\alpha_\tau} \sum_i \mathbf{1}(Y_i = h_i) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_i (1 - \underbrace{\mathbf{1}(Y_i \neq h_i)}_{\epsilon_t}) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

ϵ_t is what the classifier $h_t(\cdot)$ is trying to minimise

Let us choose α in order to minimize the global risk

$$\mathcal{R} = e^{-\alpha_\tau} - e^{-\alpha_\tau} \epsilon_t + e^{\alpha_\tau} \epsilon_t = e^{-\alpha_\tau} (1 - \epsilon_t) + e^{\alpha_\tau} \epsilon_t$$

How does one set the weight and α at each time steps?

Assume we have done the job until time τ !

$$\mathcal{R} = \sum_i \omega_i^\tau e^{-Y_i \alpha_\tau} h_\tau(\vec{x}_i)$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_{i \in \text{OK}} \omega_i^\tau + e^{\alpha_\tau} \sum_{i \in \text{NOT OK}} \omega_i^\tau = e^{-\alpha_\tau} \sum_i \mathbf{1}(Y_i = h_i) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

$$\mathcal{R} = e^{-\alpha_\tau} \sum_i (1 - \underbrace{\mathbf{1}(Y_i \neq h_i)}_{\epsilon_t}) \omega_i^\tau + e^{\alpha_\tau} \sum_i \omega_i^\tau \mathbf{1}(Y_i \neq h_i)$$

ϵ_t is what the classifier $h_t(\cdot)$ is trying to minimise

Let us choose α in order to minimize the global risk

$$\mathcal{R} = e^{-\alpha_\tau} - e^{-\alpha_\tau} \epsilon_t + e^{\alpha_\tau} \epsilon_t = e^{-\alpha_\tau} (1 - \epsilon_t) + e^{\alpha_\tau} \epsilon_t$$

$$\partial_{\alpha_t} \mathcal{R} = 0 \rightarrow \alpha_t = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}$$

ADABOOST TRAINING

$$\forall i : \omega_i^{t=1} = \frac{1}{n}$$

for $t = 1, \dots, T_{\max}$ **Do**

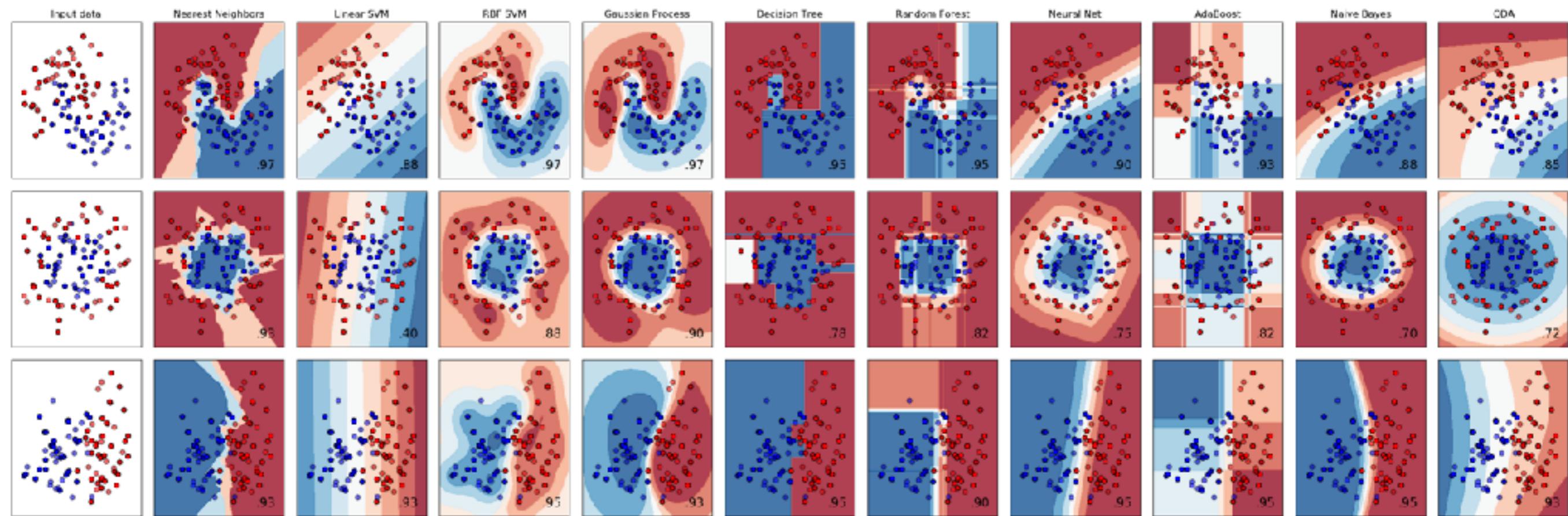
Run a classifier for : $\epsilon_t = \left[\sum_i \omega_i^\top \mathbf{1}(Y_i \neq h_i) \right]$

Set: $\alpha_t = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}$

Aggregate classifier: $H(\cdot) = H^{t-1}(\cdot) + \alpha_t h_t(\cdot)$

Update weights: $\omega_i^{t+1} = \frac{\omega_i^t e^{-Y_i \alpha_t h_t(\vec{x}_i)}}{\sum_i \omega_i^t e^{-Y_i \alpha_t h_t(\vec{x}_i)}}$

A tour on standard machine learning classifiers



From sk-learn

PCA and SVD

Blackboard lecture