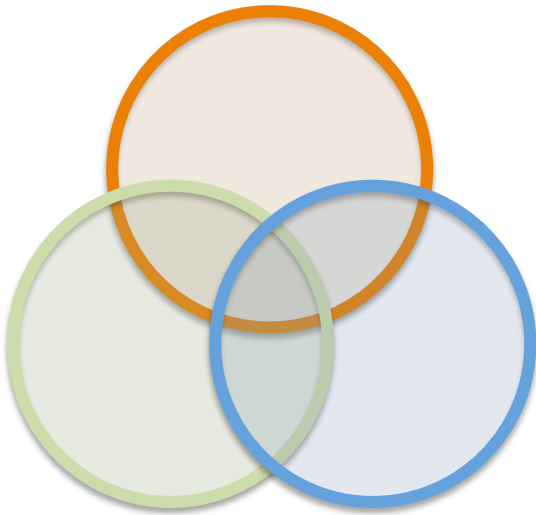


# Machine Learning for the Quantified Self



## Chapter 7

### Predictive Modeling Without Notion of Time

# Predictive modeling without notion of time

---

- Ok, let us start learning
- We will consider learning algorithms that do not take time into account explicitly
  - Of course we do have the temporal features we identified before
- We first need to consider the learning setup
  - Do we want to learn individual models?
  - Do we want to learn to predict well for the future (only for temporal data), or for unseen individuals?

# Learning setup

- Here are the options we have (validation set done in the same way):

	<i>Data</i>	
<i>Level</i>	temporal	non-temporal
individual	$\mathbf{X}_{\text{train}, \text{qs}_i} = x_{1, \text{qs}_i}, \dots, x_{n_{\text{train}}, \text{qs}_i}$ $\mathbf{X}_{\text{test}, \text{qs}_i} = x_{n_{\text{train}}+1, \text{qs}_i}, \dots, x_{N_{\text{qs}_i}, \text{qs}_i}$	$\mathbf{X}_{\text{train}, \text{qs}_i} \subset \mathbf{X}_{\text{qs}_i}$ where $ \mathbf{X}_{\text{train}, \text{qs}_i}  = n_{\text{train}} \wedge \mathbf{X}_{\text{train}, \text{qs}_i} \cap \mathbf{X}_{\text{test}, \text{qs}_i} = \emptyset$ $\mathbf{X}_{\text{test}, \text{qs}_i} \subset \mathbf{X}_{\text{qs}_i}$ where $ \mathbf{X}_{\text{test}, \text{qs}_i}  = (N_{\text{qs}_i} - n_{\text{train}}) \wedge \mathbf{X}_{\text{train}, \text{qs}_i} \cap \mathbf{X}_{\text{test}, \text{qs}_i} = \emptyset$
population -	$\mathbf{X}_{\text{train}} \subset \{\mathbf{X}_{\text{qs}_1}, \dots, \mathbf{X}_{\text{qs}_n}\}$ where $ \mathbf{X}_{\text{train}}  = n_{\text{train}} \wedge \mathbf{X}_{\text{train}} \cap \mathbf{X}_{\text{test}} = \emptyset$	
unknown users	$\mathbf{X}_{\text{test}} \subset \{\mathbf{X}_{\text{qs}_1}, \dots, \mathbf{X}_{\text{qs}_n}\}$ where $ \mathbf{X}_{\text{test}}  = (n - n_{\text{train}}) \wedge \mathbf{X}_{\text{train}} \cap \mathbf{X}_{\text{test}} = \emptyset$	
population -	$\mathbf{X}_{\text{train}} = \{\mathbf{X}_{\text{train}, \text{qs}_1}, \dots, \mathbf{X}_{\text{train}, \text{qs}_n}\}$ (Note that we refer to the definition of the training sets specified in the row for the individual level here)	
unseen data of known users	$\mathbf{X}_{\text{test}} = \{\mathbf{X}_{\text{test}, \text{qs}_1}, \dots, \mathbf{X}_{\text{test}, \text{qs}_n}\}$	

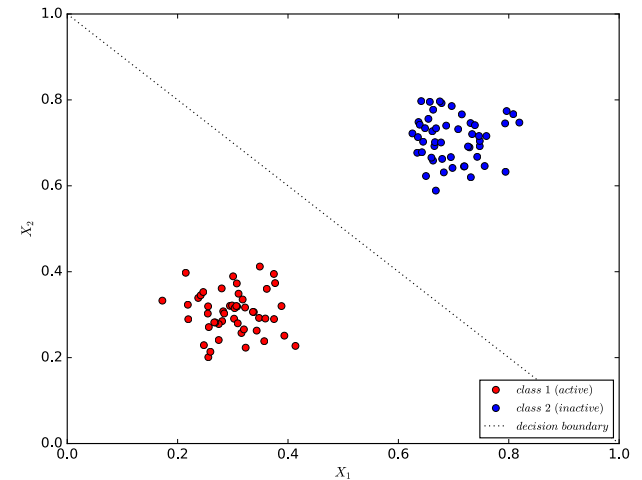
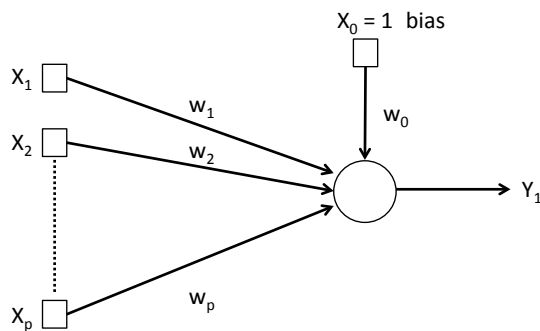
# Learning algorithms

---

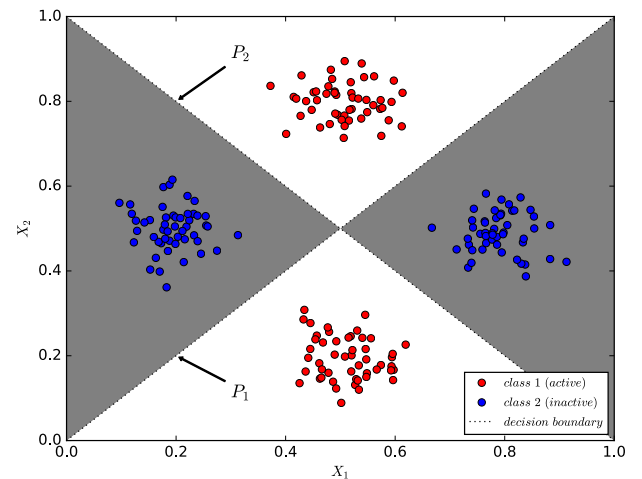
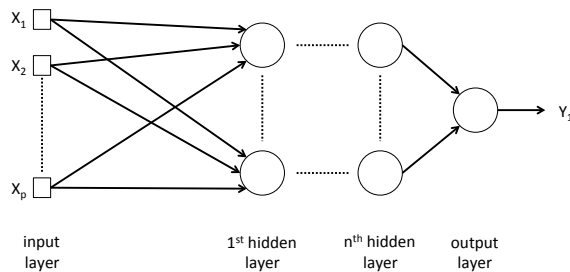
- Assume you know the basic learning algorithms
- Quick recap, two slides per learning algorithm
  - Neural Networks
  - Support Vector Machines
  - K-Nearest Neighbor
  - Decision trees
  - Naïve Bayes
  - Ensembles

# Neural Networks (1)

- Perceptron



- Multi-Layer Perceptron

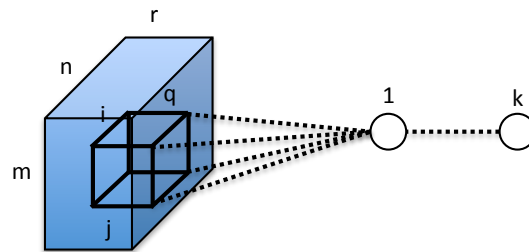


# Neural Networks (2)

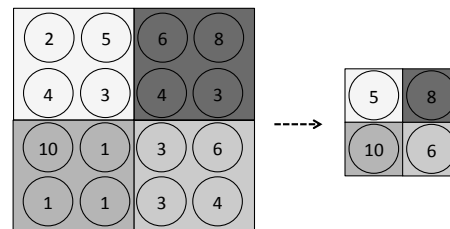
- Convolutional neural networks

- Extract features using:

- Convolutional layers



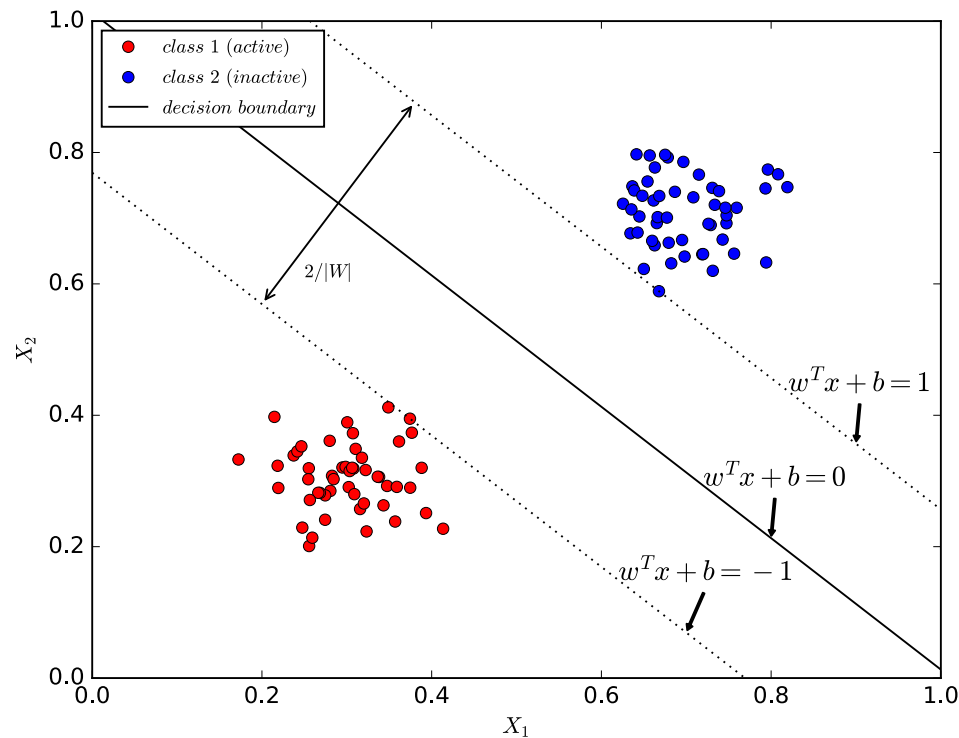
- Pooling layers



- Followed by a “regular” neural network

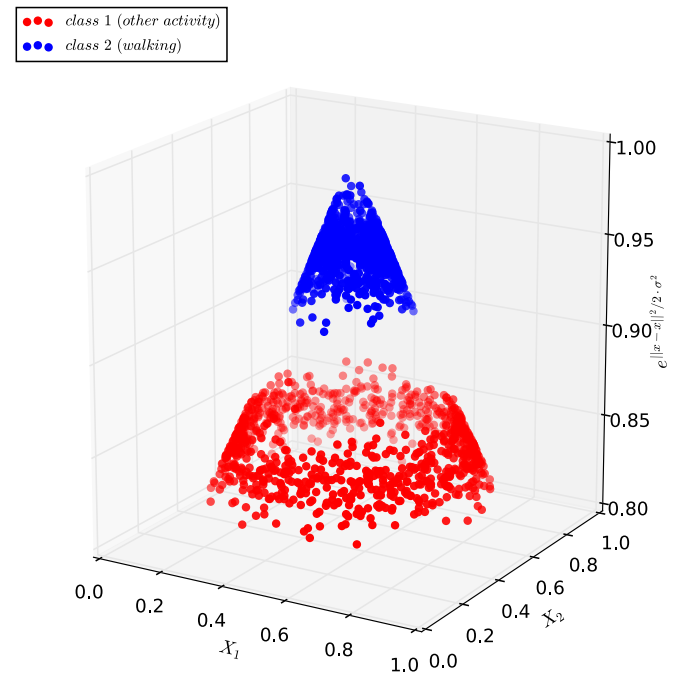
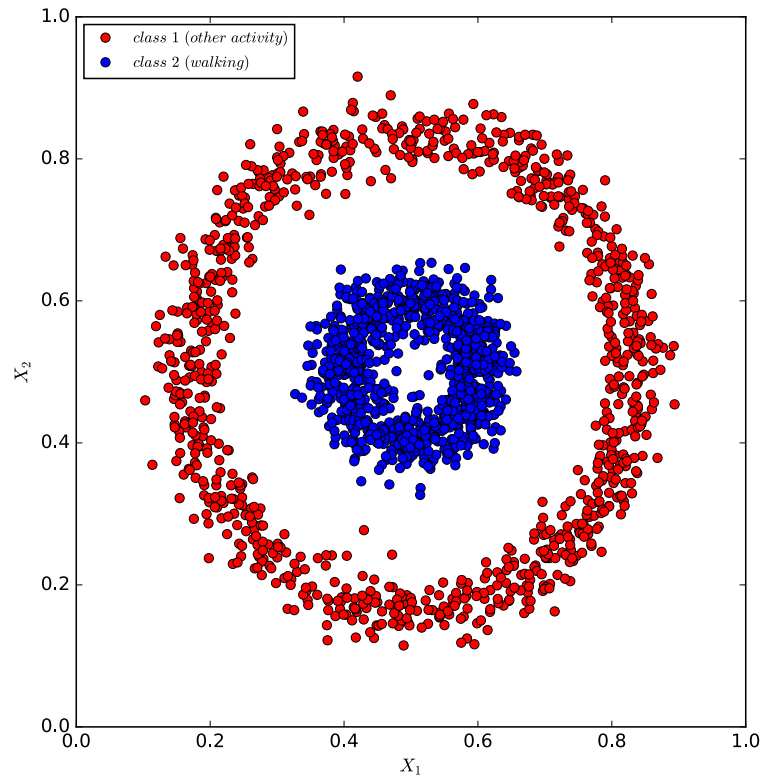
# Support Vector Machines (1)

- Find hyperplane that maximizes separation between classes



# Support Vector Machines (2)

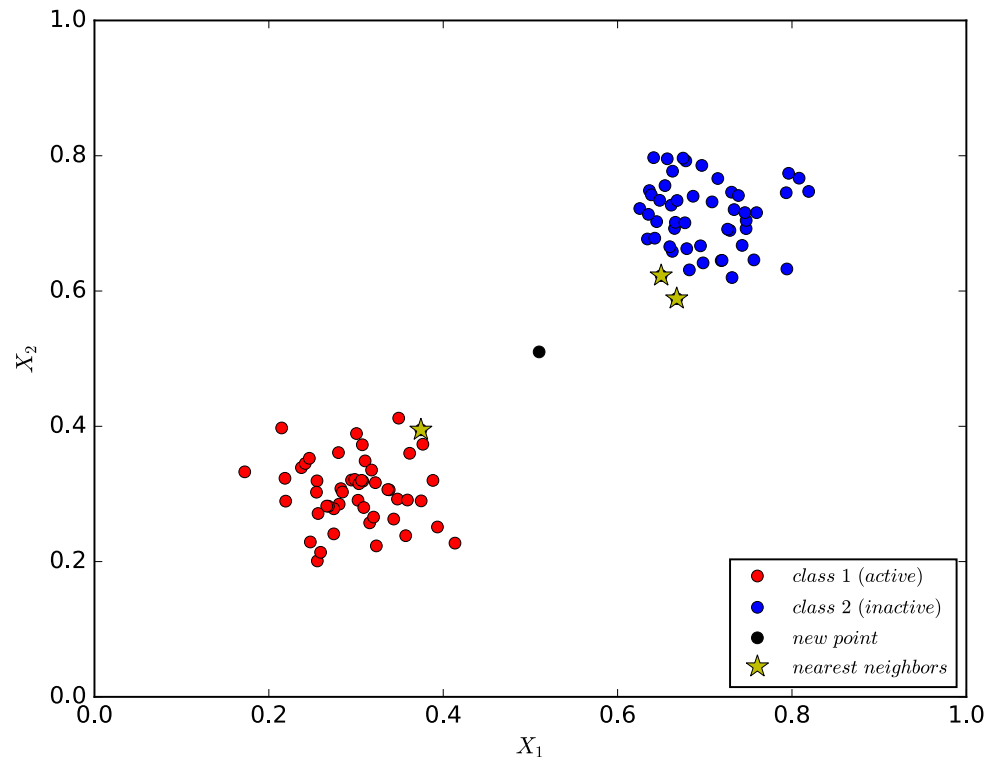
- Kernel functions





# K-Nearest Neighbor (1)

- Find  $k$  closest examples
  - use distance functions discussed before



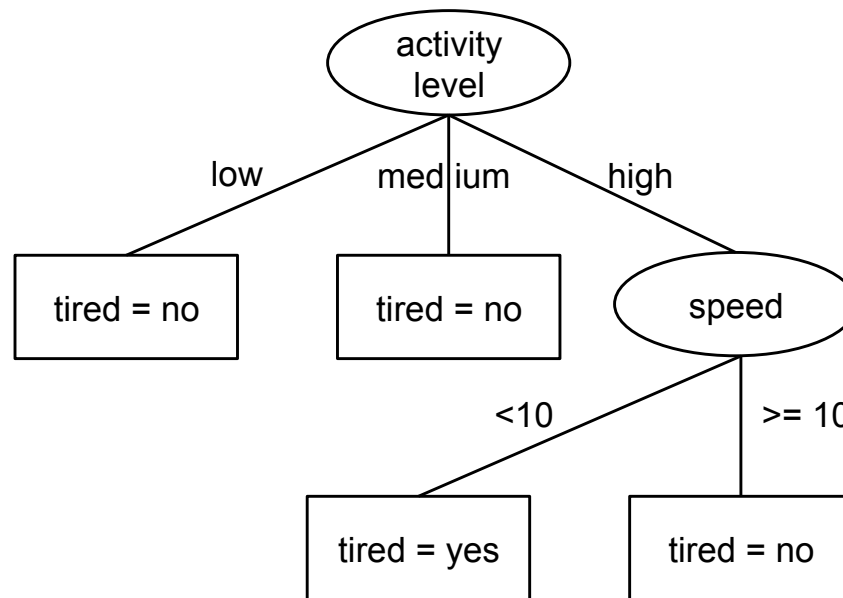
# K-Nearest Neighbor (2)

---

- Assign classes based on some function, for example:
  - Majority class
  - Average value

# Decision tree (1)

- Create a tree to decide on target value



# Decision tree (2)

- Use criterion to decide on which attributes are most important (i.e. should be selected first)

– Categorical target:

$$entropy(\mathbf{G}) = \sum_{v \in \mathcal{G}} -p_v \log_2 p_v \text{ where } p_v = \frac{|\{g_i \in \mathbf{G} | g_i = v\}|}{|\mathbf{G}|}$$
$$gain(X_i, \mathbf{X}, \mathbf{G}) = entropy(\mathbf{G}) - \sum_{v \in \mathcal{X}_i} \frac{|\{x_j \in \mathbf{X} | x_j^i = v\}|}{|\mathbf{X}|} \cdot entropy(\{g_i \in \mathbf{G} | x_j^i = v\})$$

– Numerical target:

$$sd\_red(X_i, \mathbf{X}, \mathbf{Y}) = \sigma_Y - \sum_{v \in \mathcal{X}_i} \frac{|\{y_i \in \mathbf{Y} | x_j^i = v\}|}{|\mathbf{Y}|} \sigma_{\{y_i \in \mathbf{Y} | x_j^i = v\}}$$

# Naïve Bayes

- We use Bayes formula to express the probability of a target value  $g$  given our observations  $x$ :

$$p(G = g|x) = \frac{p(x|G = g)p(G = g)}{p(x)}$$

# Ensembles (1)

- Bagging
  - Create multiple models on samples of the data

---

**Algorithm 14:** Bagging

---

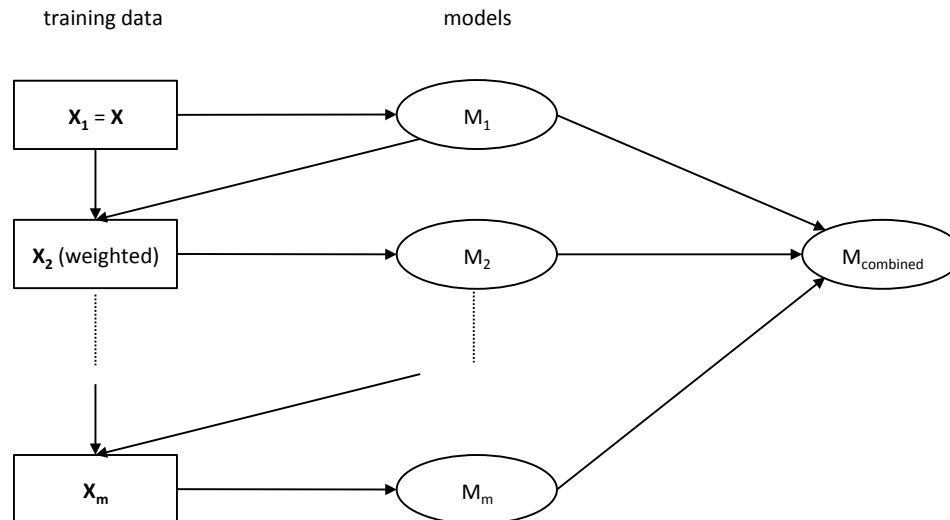
```
models = []  
for  $i$  in  $1, \dots, m$  do  
    Take a sample  $\mathbf{X}_i \subset \mathbf{X}$   
    Build a model  $M_i$  on training data  $\mathbf{X}_i$  using some learning algorithm  
    Add  $M_i$  to models  
end
```

---

- Combine output of models using some aggregation function
- Tackles *variance* problem

# Ensembles (2)

- Boosting
  - Build models sequentially and focus models on where we made mistakes before



- Focuses on *bias* problem

# Practical considerations

---

- Ok, let us start to apply the learning algorithms!
  - We need to be careful not to overfit (think of the first set of slides)
  - We can select features using forward and backward selection
  - We can consider regularization (punishing more complex models)



# Forward selection

- We iteratively add the most predictive feature:

---

**Algorithm 17:** Forward selection

---

```
selected_attributes = {}
performances = []
for  $k = 1, \dots, p$  do
    best_attribute = ""
    best_performance =  $\infty$ 
    available_attributes =  $X \setminus \text{selected\_attributes}$ 
    for  $l \in 1, \dots, |\text{available\_attributes}|$  do
        temp_attributes =  $\text{selected\_attributes} \cup \text{available\_attributes}_l$ 
        performance =  $\text{learn\_model}(\text{temp\_attributes}, \mathbf{X})$ 
        if  $\text{performance} < \text{best\_performance}$  then
            best_attribute =  $\text{available\_attributes}_l$ 
            best_performance = performance
        end
    end
    performances[k] = best_performance
    selected_attributes =  $\text{selected\_attributes} \cup \text{best\_attribute}$ 
end
return performances
```

---

# Backward selection

- We iteratively remove the least predictive feature:

---

**Algorithm 18:** Backward selection

---

```
selected_attributes = X
performances = []
performances[p] = learn_model(selected_attributes, X)
for  $k = p - 1, \dots, 1$  do
    worst_attribute = ""
    best_performance =  $\infty$ 
    available_attributes = selected_attributes
    for  $l \in 1, \dots, |available\_attributes|$  do
        temp_attributes = selected_attributes \ available_attributes $l$ 
        performance = learn_model(temp_attributes, X)
        if  $performance < best\_performance$  then
            worst_attribute = available_attributes $l$ 
            best_performance = performance
        end
    end
    performances[k] = best_performance
    selected_attributes = selected_attributes \ worst_attribute
end
return performances
```

---

# Regularization

- We add a term to the error function to punish for more complex models
  - For neural networks:

$$E = \sum_{i=0}^N (y_i - \hat{y}_i)^2 + \frac{\lambda}{2} ||w||$$

# CrowdSignals

---

- Remember, we have two tasks for our CrowdSignals dataset:
  1. Predict label (activity)
  2. Predict heart rate
- We consider the first with a non-temporal setup and the second with a temporal setup

# Predicting the label

---

- We currently have a set of binary attributes for each label value
  - We need to transform these to a single categorical attribute
  - How do we do this? What if we have multiple “1”s per instance?
  - We only use instances where we have a single activity marked with a 1
  - Results in 1837 instances (out of 2895)

# Splitting the data

---

- How do we split the data into training and test?
  - We assign 70% of the instances to the training set and 30% to the test set
  - We do so in a *stratified* way
  - Once again, we ignore the time component

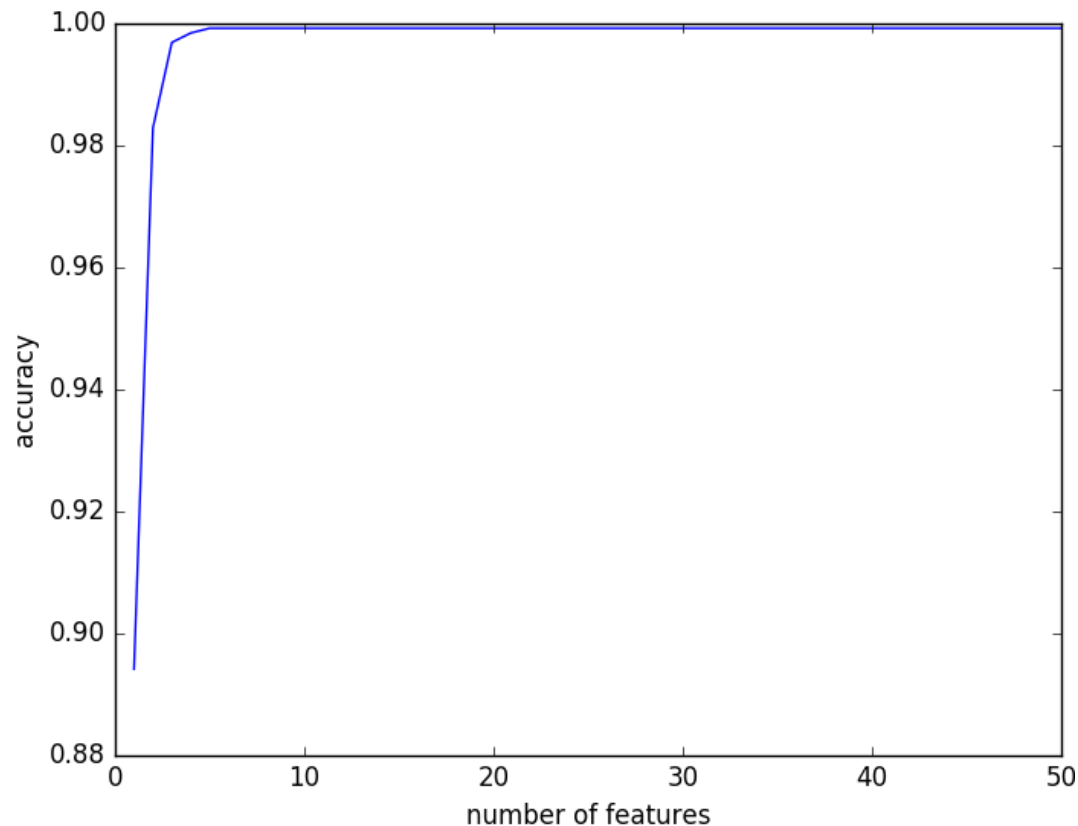
# Feature selection (1)

---

- We have 517 features if we use all the ones we have identified
  - Initial features (21)
  - PCA (7)
  - Time domain features, mean + standard deviation ( $2 \times 28$  features = 56)
  - Frequency domain (18 frequencies, 24 initial features = 432)
  - Cluster (1)
- Should we select some features?

# Feature selection (2)

- Forward selection:





# Feature selection (3)

- We end up with different options which we will try:

Dataset name	Basic features	PCA	Time-and frequency features	fre-based	Clusters	#features
Initial set (cleaned)	21					21
Chapter 3	21	7				28
Chapter 4	21	7	56 + 432			516
Chapter 5	21	7	56 + 432	1		517
Selected features			2 + 8			10

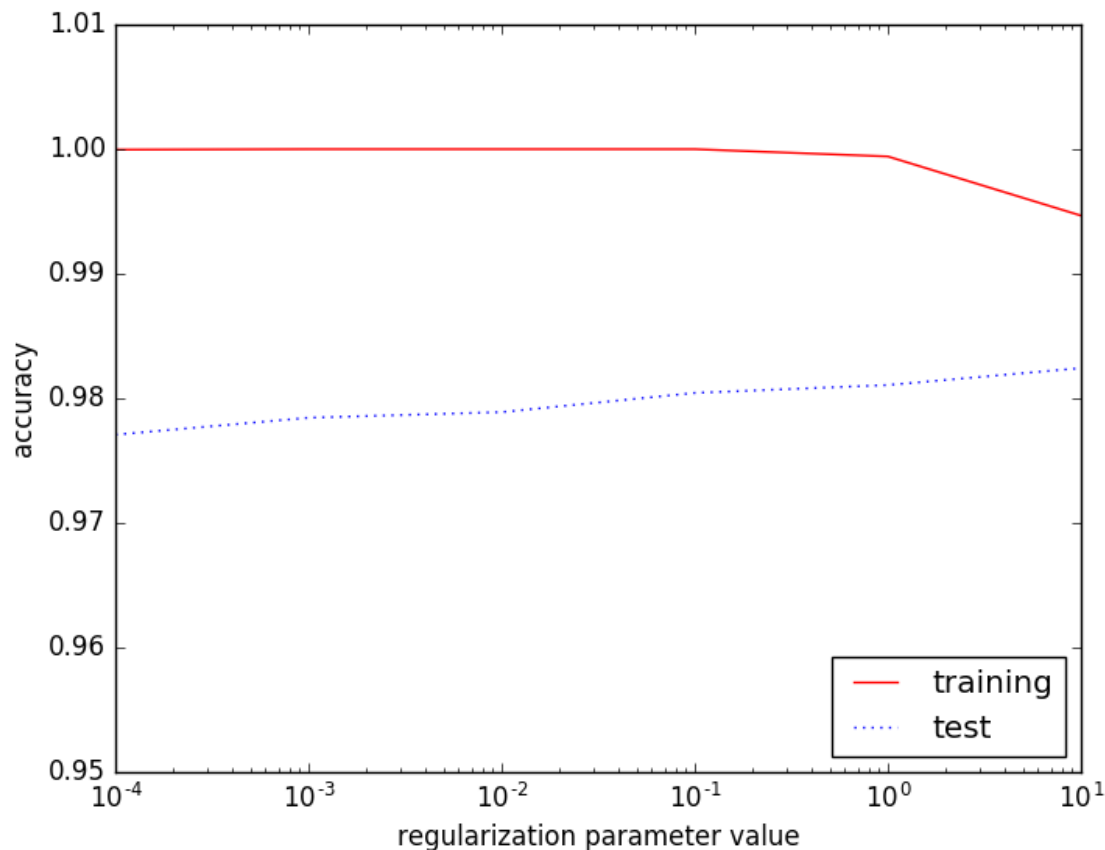
# Influence of parameter settings (1)

---

- How does model complexity influence our performance in this case?
  - Remember the theoretical results we have seen
  - We can tailor the parameters of the models
  - We can change the weight of the regularization terms

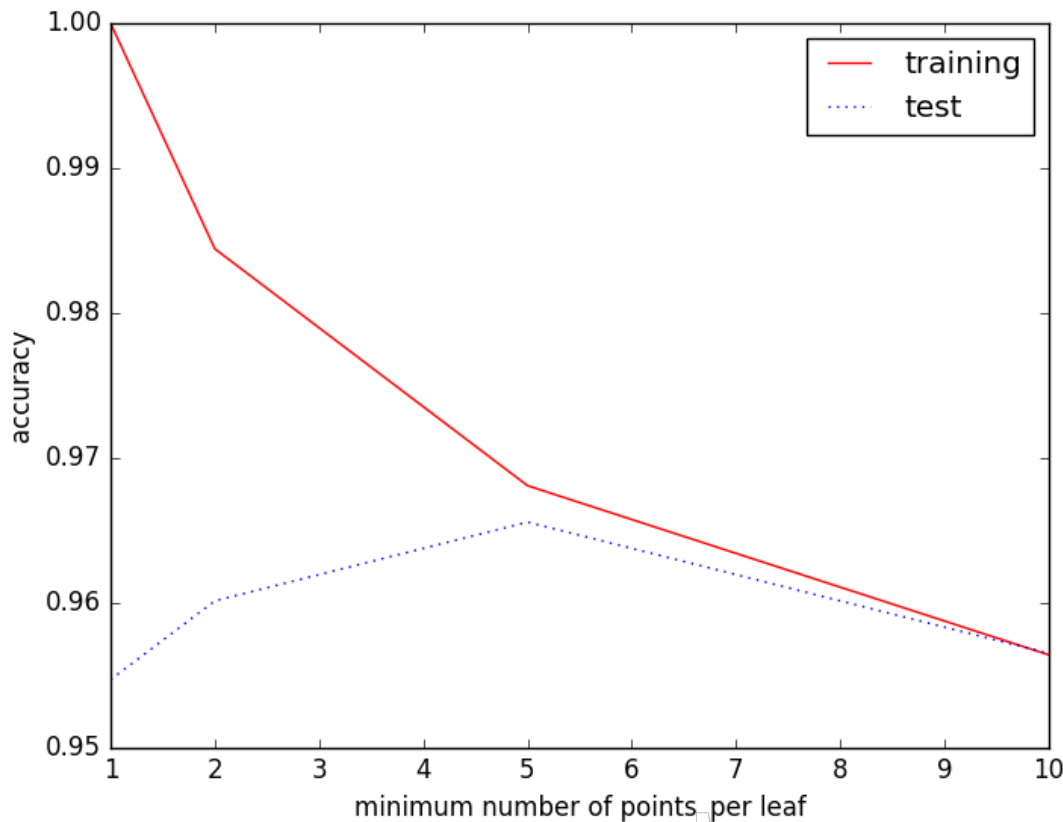
# Influence of parameter settings (2)

- Neural network and regularization term



# Influence of parameter settings (3)

- Decision tree and points per leaf:

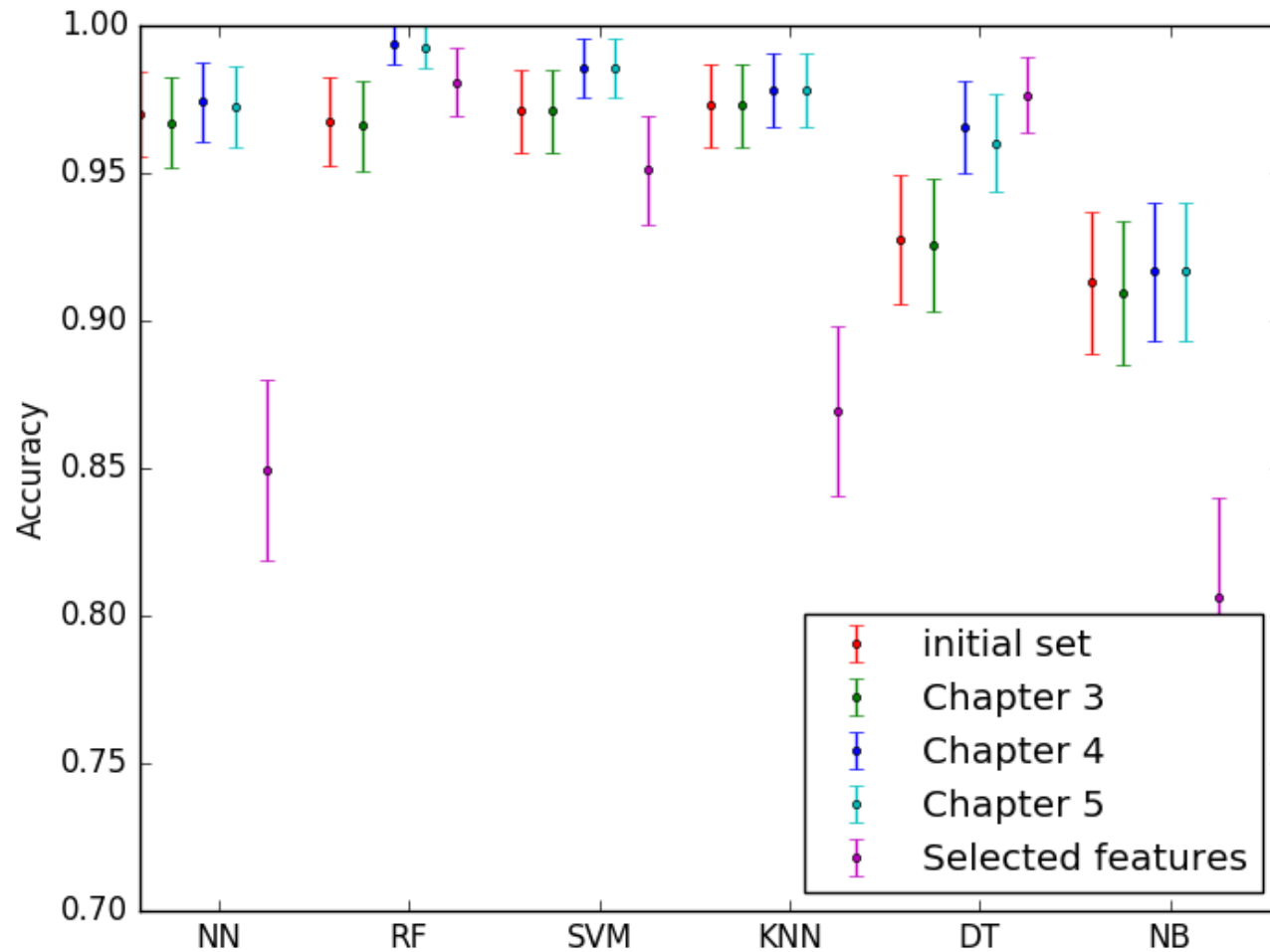


# Influence of parameter settings (4)

- We try different parameter settings
  - study the generalizability based on a cross validation on the training set

Algorithm	Variant description	Parameters varied
Neural Network (NN)	Multi-Layer Perceptron with 1 or 2 hidden layers, a logistic activation function and one output node per class	hidden layer composition: {single layer with 5, 10, 25, or 100 neurons, two layer with 100 and 5 neurons or 100 and 10} maximum iterations: {1000, 2000}
Support Vector Machines (SVM)	SVM with kernel function with one SVM model per class	maximum iterations: {1000, 2000} C: {1, 10, 100} tolerance = {0.001, 0.0001} kernel function: { <i>rbf</i> , <i>polynomial</i> }
K-Nearest Neighbor (KNN)	KNN model using simple Euclidean distance	k: {1, 2, 5, 10}
Decision Tree (DT)	Decision tree algorithm following the CART approach	minimum samples per leaf: {2, 10, 50, 100, 200} splitting criterion: { <i>gini</i> , <i>entropy</i> }
Naive Bayes (NB)	Basic Naive Bayes approach	-
Random Forest (RF)	Basic Random forest approach	minimum samples per leaf: {2, 10, 50, 100, 200} number of trees: {10, 50, 100} splitting criterion: { <i>gini</i> , <i>entropy</i> }

# Results test set (1)



# Results test set (2)

- How did we computer the 95% confidence intervals?
  - $a$  is the accuracy
  - $n$  the number of samples in the test set

$$sd = \sqrt{a(1 - a)/n}$$

# Results test set (3)

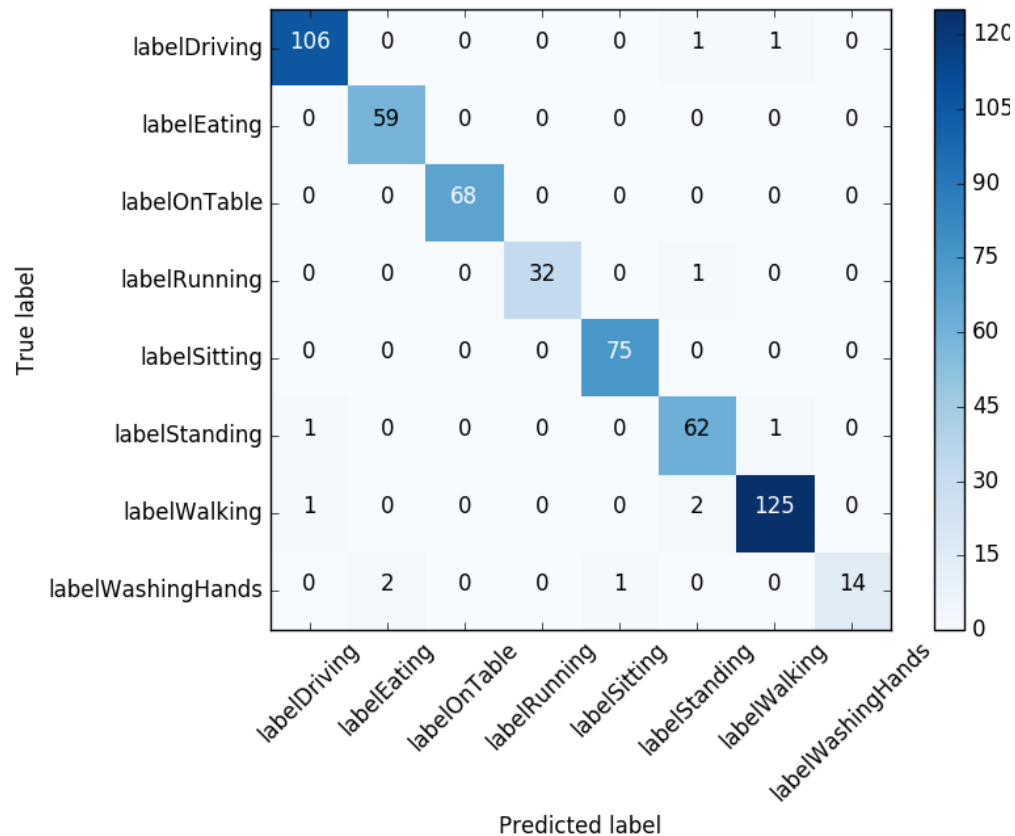
- Important features (random forest):

Feature	Importance
<i>gyr_phone_x_temp_std_ws_120</i>	0.2855
<i>press_phone_pressure_temp_mean_ws_120</i>	0.2499
<i>acc_phone_y_freq_0.0_Hz_ws_40</i>	0.2433
<i>mag_watch_y_pse</i>	0.0873
<i>mag_phone_z_max_freq</i>	0.0341
<i>gyr_watch_y_freq_weighted</i>	0.0284
<i>gyr_phone_y_freq_1.0_Hz_ws_40</i>	0.0255
<i>acc_phone_x_freq_1.9_Hz_ws_40</i>	0.0227
<i>acc_watch_y_freq_0.5_Hz_ws_40</i>	0.0164
<i>mag_watch_z_freq_0.9_Hz_ws_40</i>	0.0068



# Results test set (4)

- Confusion matrix:



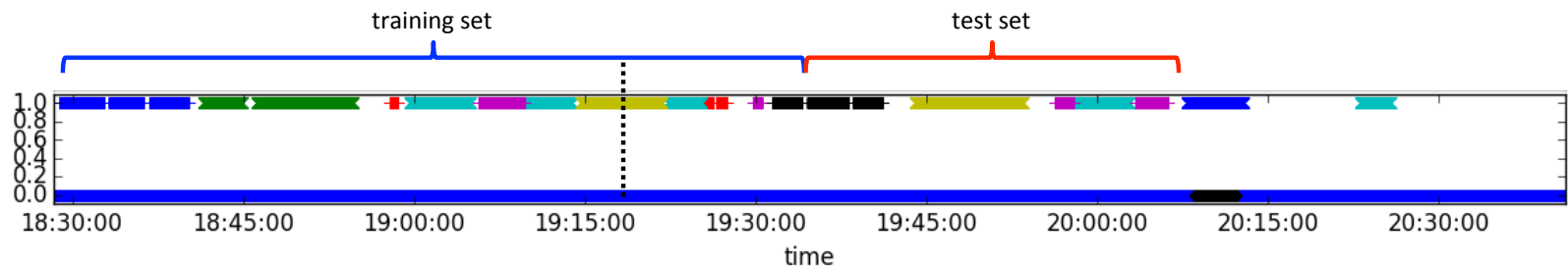
# Results test set (5)

---

- Why are the results this good?
  - Dataset has limited variation
  - Spans only short time period
  - ...

# Predicting the heart rate

- We take it as a temporal problem
  - Want to learn to predict future time points well
- How do we split the training and test set here?



- Why this way?

# Features

Dataset name	Basic features	PCA	Temporal features	Clusters	#features
Initial set (cleaned)	21				21
Chapter 3	28	7			35
Chapter 4	28	7	71 + 432		538
Chapter 5	28	7	71 + 432	1	539
Selected features	1	1	6 + 2		10

# Algorithms we tried

Algorithm	Variant description	Parameters varied
Neural Network (NN)	Mutli-Layer Perceptron with 1 or 2 hidden layers and an identity activation function and one output node per class	hidden layer composition: {single layer with 5, 10, 25, or 100 neurons, two layer with 100 and 5 neurons or 100 and 10} maximum iterations: {1000, 2000}
Support Vector Regression (SVR)	Linear SVR (no kernel function)	maximum iterations: {1000, 2000} C: {1, 10, 100} tolerance = {0.001, 0.0001}
K-Nearest Neighbor (KNN)	KNN model using simple euclidean distance	k: {1, 2, 5, 10}
Decision Tree (DT)	Decision tree algorithm following the CART approach with the mean squared error as a splitting criterion	minimum samples per leaf: {50, 100, 200}
Random Forest (RF)	Basic Random forest approach with the mean squared error as a splitting criterion	minimum samples per leaf: {50, 100, 200} number of trees: {10, 50, 100}

# Results (1)

Approach	NN		RF		SVM		KNN		DT	
Features	training	test	training	test	training	test	training	test	training	test
initial set	725.3 (1009.3)	1584.6 (1417.6)	582.5 (857.1)	1548.0 (1346.5)	1155.8 (1216.8)	1303.3 (1067.2)	186.8 (574.1)	2309.7 (2027.9)	542.8 (855.8)	1496.7 (1268.4)
Chapter 3	716.3 (971.6)	1510.0 (1333.6)	476.0 (800.9)	1750.4 (1500.1)	1600.2 (4785.8)	1465.2 (1609.4)	186.8 (574.1)	2309.7 (2027.9)	537.2 (850.8)	1501.4 (1282.2)
Chapter 4	523.6 (846.2)	1968.0 (2801.3)	15.1 (41.6)	1310.3 (1129.2)	1516.1 (2783.7)	1967.4 (3950.7)	134.2 (473.6)	2259.0 (1848.2)	368.2 (643.8)	1317.0 (826.8)
Chapter 5	557.7 (861.5)	2014.7 (2114.9)	7.0 (21.1)	1384.8 (1292.9)	1211.8 (2051.6)	1997.2 (3243.5)	134.2 (473.6)	2259.0 (1848.2)	368.2 (643.8)	1317.0 (826.8)
Selected features	768.9 (1369.2)	2825.3 (2855.4)	620.7 (1109.3)	2614.8 (2613.4)	808.5 (1909.1)	2598.2 (2944.2)	286.1 (719.9)	2571.1 (2896.0)	538.6 (1039.8)	2740.0 (2878.4)

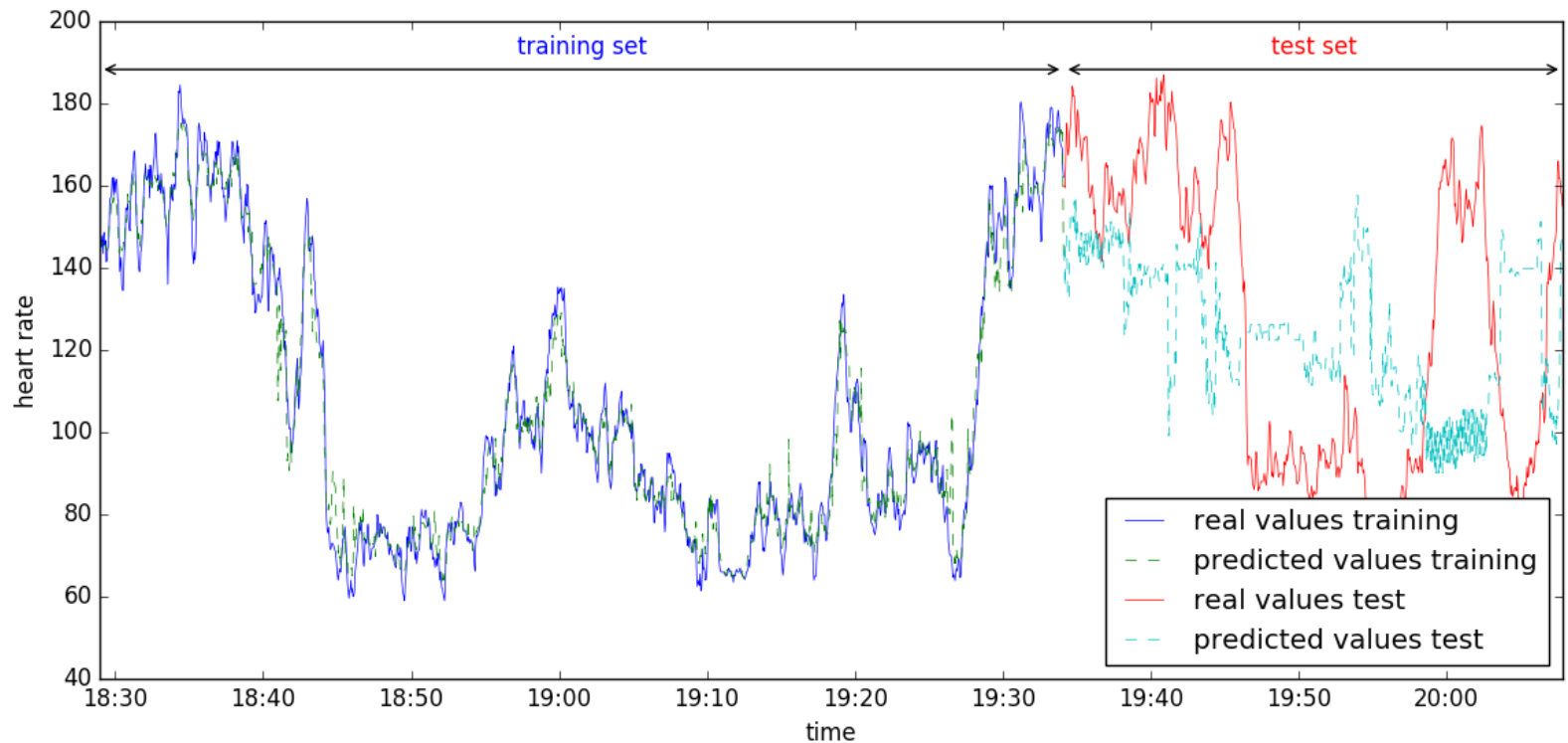
# Results (2)

- Standard deviation computed based on predictions per instance on the test set
- Predictors random forest:

Feature	Importance
<i>acc_watch_y_temp_mean_ws_120</i>	0.4631
<i>press_phone_pressure_temp_mean_ws_120</i>	0.1398
<i>mag_watch_x_temp_mean_ws_120</i>	0.0465
<i>temp_pattern_labelEating(b)labelEating</i>	0.0404
<i>mag_watch_y_temp_mean_ws_120</i>	0.0385
<i>temp_pattern_labelEating</i>	0.0383
<i>labelEating</i>	0.0382
<i>pca_4_temp_mean_ws_120</i>	0.0208
<i>pca_7</i>	0.0131
<i>pca_7_temp_mean_ws_120</i>	0.0112

# Results (3)

- Visual predictions:





# Why is it so bad?

- Ideas?

