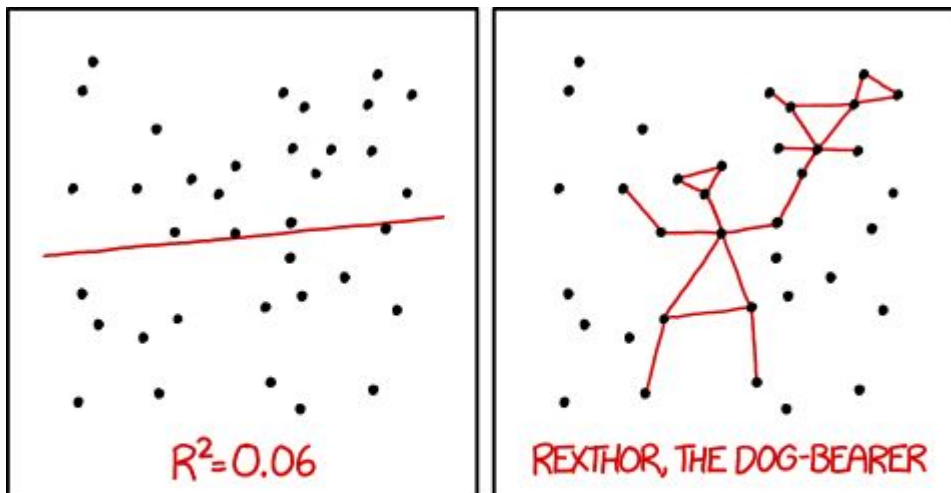


# ADVANCED REGRESSION

**Probabilistic interpretation of LR. Classification  
algorithms for regression**

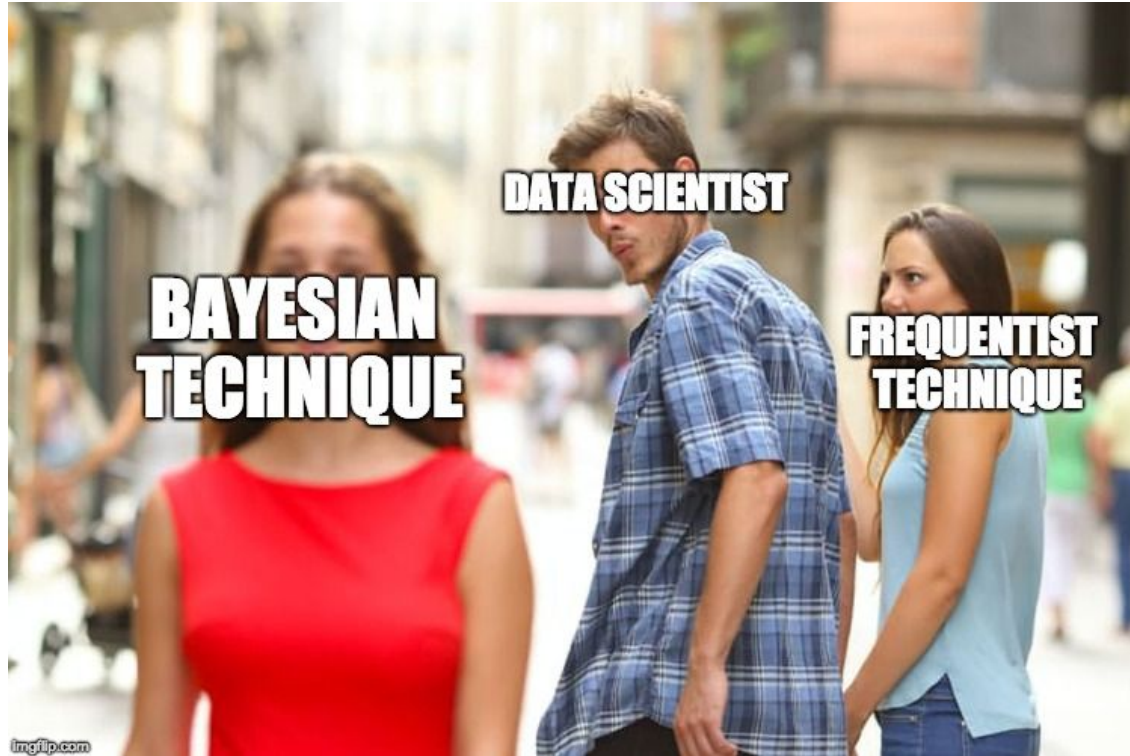
# CONTENTS

1. Bayesian explanation of regularized regression
2. Classification algorithms for regression



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER  
TO GUESS THE DIRECTION OF THE CORRELATION FROM THE  
SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

# BAYESIAN REGRESSION



# CLASSICAL PROBABILISTIC VIEW ON LINEAR REGRESSION

Consider we have  $n$  points  $\mathbf{Y}$  drawn i.i.d. from the normal distribution. The probability of those points being drawn defines the likelihood function which is just a multiplication of their densities in every points.

$$\mathcal{L}(\mu|y) = \prod_{i=1}^n P_Y(y_i|\mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y_i-\mu)^2}{2\sigma^2}} \quad (1)$$

A good estimate of mean maximizes that likelihood

$$\begin{aligned} \hat{\mu} &= \arg \max_{\mu} \mathcal{L}(\mu|y) = \arg \max_{\mu} \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y_i-\mu)^2}{2\sigma^2}} \\ &= \arg \max_{\mu} \log \left( \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y_i-\mu)^2}{2\sigma^2}} \right) \\ &= \arg \max_{\mu} \sum_{i=1}^n \log \left( \frac{1}{\sigma\sqrt{2\pi}} \right) + \log \left( e^{-\frac{(y_i-\mu)^2}{2\sigma^2}} \right) \\ &= \arg \max_{\mu} \sum_{i=1}^n \log \left( e^{-\frac{(y_i-\mu)^2}{2\sigma^2}} \right) \\ &= \arg \max_{\mu} \sum_{i=1}^n -\frac{(y_i - \mu)^2}{2\sigma^2} \\ &= \arg \min_{\mu} \sum_{i=1}^n (y_i - \mu)^2 \quad (2) \end{aligned}$$

# CLASSICAL PROBABILISTIC VIEW ON LINEAR REGRESSION

- Assume our mean is a function of predictors  $\mathbf{x}$

$$E(y|\mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

- Thus our target is distributed according to

$$y \sim N(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2)$$

- Estimating regression parameters given (2)

$$\begin{aligned}\beta &= \arg \min_{\beta} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p))^2 \\ &= \arg \min_{\beta} \sum_{i=1}^n (y_i - \hat{y}_i)^2\end{aligned}$$

# A PROBABILISTIC INTERPRETATION OF REGULARIZATION

- Using a Bayes theorem we can estimate the **probability distribution** of the parameters  $\theta$  given the data we observe  $\mathbf{Y}$

$$P(\theta|y) = \frac{P(y|\theta)P(\theta)}{P(y)}$$
$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}}$$

- That gives us an opportunity to set the a prior distribution of model parameters
- Compare that with the classical method where we instead try to find the best parameters to maximize the **likelihood of parameters** given the data

# A PROBABILISTIC INTERPRETATION OF REGULARIZATION

- We maximize the posterior probability estimate using Bayes theorem (Maximum A Posteriori estimation)

$$\begin{aligned}\hat{\theta}_{\text{MAP}} &= \arg \max_{\theta} P(\theta|y) \\ &= \arg \max_{\theta} \frac{P(y|\theta)P(\theta)}{P(y)} \\ &= \arg \max_{\theta} P(y|\theta)P(\theta) \\ &= \arg \max_{\theta} \log(P(y|\theta)P(\theta)) \\ &= \arg \max_{\theta} \log P(y|\theta) + \log P(\theta)\end{aligned}$$

- Compare that to MLE estimate

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \log P(y|\theta)$$

# A PROBABILISTIC INTERPRETATION OF L2 REGULARIZATION

- Assume our model parameters zero-mean normally distributed with  $\tau^2$  variance (prior knowledge)

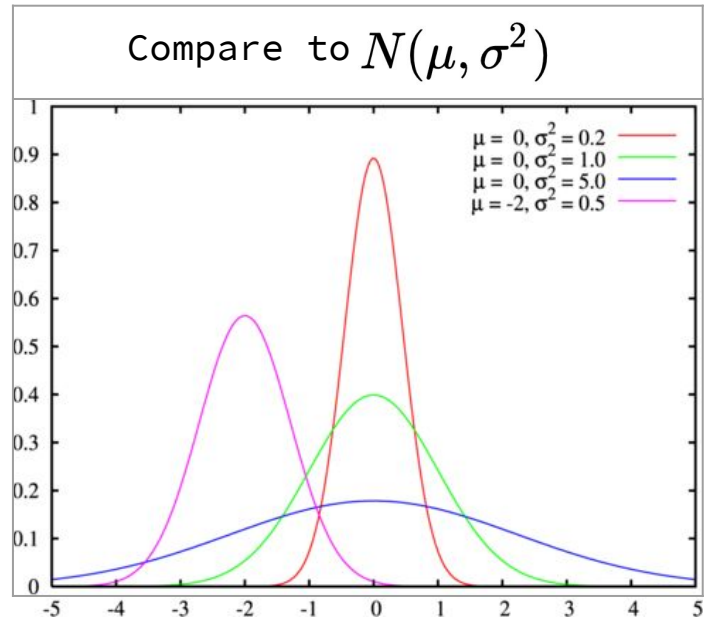
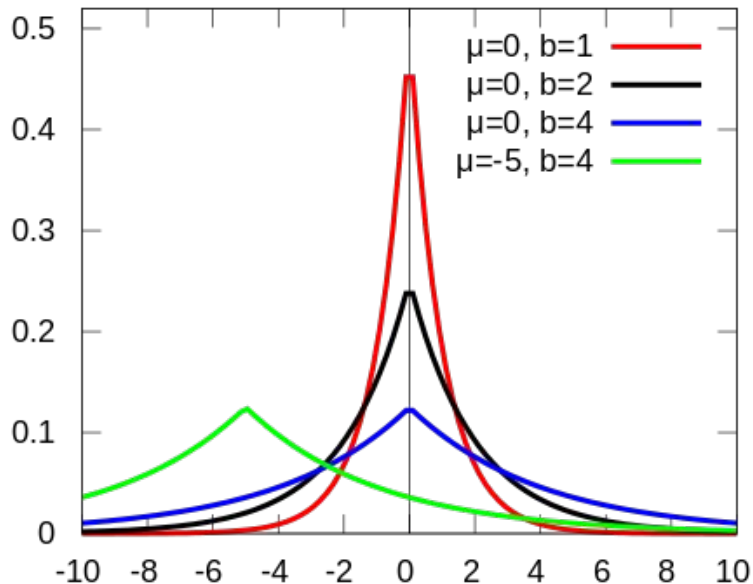
$$\begin{aligned} & \arg \max_{\beta} \left[ \log \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}))^2}{2\sigma^2}} + \log \prod_{j=0}^p \frac{1}{\tau \sqrt{2\pi}} e^{-\frac{\beta_j^2}{2\tau^2}} \right] \\ &= \arg \max_{\beta} \left[ - \sum_{i=1}^n \frac{(y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}))^2}{2\sigma^2} - \sum_{j=0}^p \frac{\beta_j^2}{2\tau^2} \right] \\ &= \arg \min_{\beta} \frac{1}{2\sigma^2} \left[ \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}))^2 + \frac{\sigma^2}{\tau^2} \sum_{j=0}^p \beta_j^2 \right] \\ &= \arg \min_{\beta} \left[ \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}))^2 + \lambda \sum_{j=0}^p \beta_j^2 \right] \end{aligned}$$

- Small variance  $\tau^2$  (large  $\lambda$ ) leads to reduction of coefficients. If we have a large variance (small  $\lambda$ ) the coefficients are not affected much.



# A PROBABILISTIC INTERPRETATION OF L1 REGULARIZATION

- Laplace distribution with mean  $\mu$  and diversity  $b$  defined by a probability density function  $Laplace(\mu, b) = \frac{1}{2b} e^{-\frac{|x-\mu|}{b}}$



# A PROBABILISTIC INTERPRETATION OF L1 REGULARIZATION

- Assume our model parameters zero-mean Laplace-distributed with diversity **b**

$$\begin{aligned} & \arg \max_{\beta} \left[ \log \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}))^2}{2\sigma^2}} + \log \prod_{j=0}^p \frac{1}{2b} e^{-\frac{|\beta_j|}{b}} \right] \\ &= \arg \max_{\beta} \left[ - \sum_{i=1}^n \frac{(y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}))^2}{2\sigma^2} - \sum_{j=0}^p \frac{|\beta_j|}{b} \right] \\ &= \arg \min_{\beta} \frac{1}{2\sigma^2} \left[ \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}))^2 + \frac{2\sigma^2}{b} \sum_{j=0}^p |\beta_j| \right] \\ &= \arg \min_{\beta} \left[ \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}))^2 + \lambda \sum_{j=0}^p |\beta_j| \right] \end{aligned}$$

- L1 regularization promotes sparsity in comparison with “just reducing the coefficients” in L2. That makes sense if you look at Laplacean density where there is a sharp increase in  $x = \text{mean}$ .

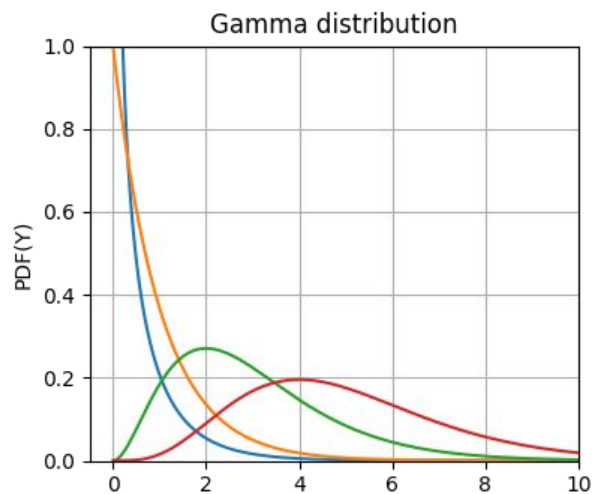
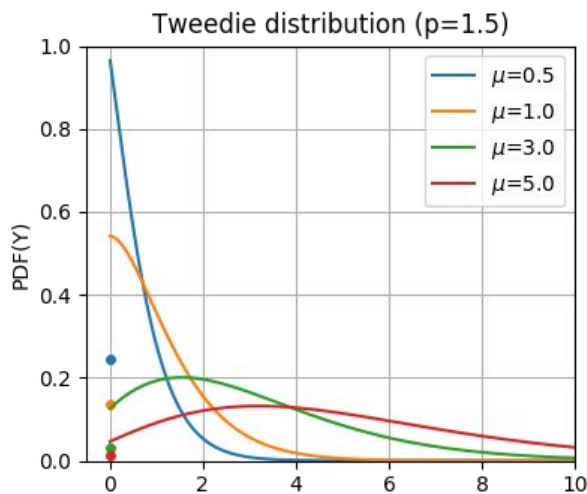
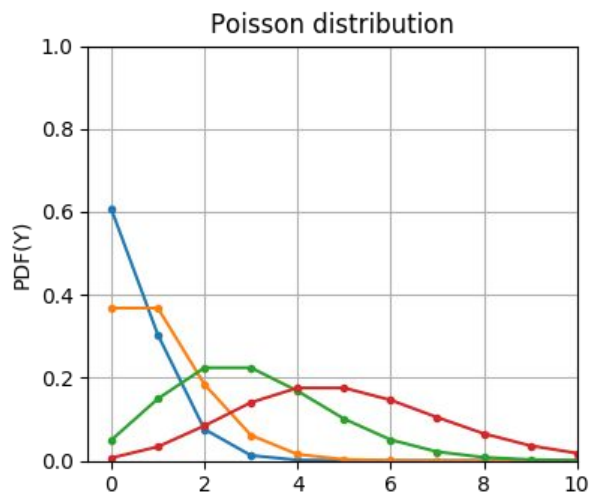
WHY L1 ZERO OUT COEFFICIENTS WHEREAS L2 DOES NOT?

# WHY L1 ZERO OUT COEFFICIENTS WHEREAS L2 DOES NOT?

- Laplace distribution (sharp in  $x=\text{mean}$ ) vs Normal distribution (smooth)
- Intuitive understanding through Gradient Descent  
<https://developers.google.com/machine-learning/crash-course/regularization-for-sparsity/l1-regularization>
- Intuitive understanding through visualization in 2d case  
<https://explained.ai/regularization/L1vsL2.html>

# GENERALIZED LINEAR MODELS

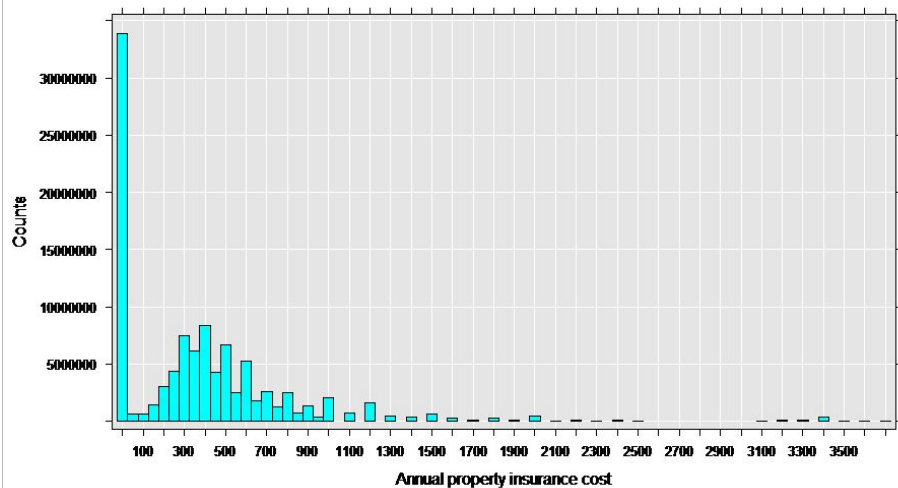
- What if we change the hypothesis  $y \sim N(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2)$



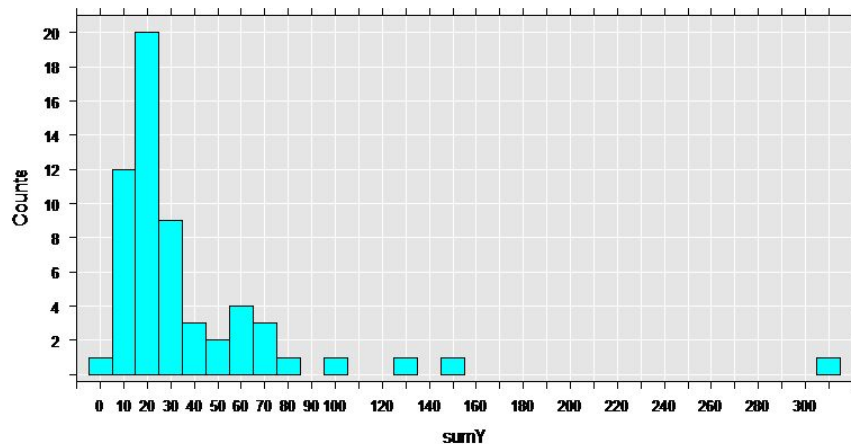
- [https://scikit-learn.org/stable/modules/linear\\_model.html#generalized-linear-regression](https://scikit-learn.org/stable/modules/linear_model.html#generalized-linear-regression)

# REAL WORLD EXAMPLES

Insurance cost  
(Tweedie distribution)



Number of calls arriving in a call  
center per hour  
(Poisson distribution)



# CLASSIFICATION ALGORITHMS FOR REGRESSION



# KNN REGRESSOR

How to calculate continuous variable for KNN?



# KNN REGRESSOR

How to calculate continuous variable for KNN?

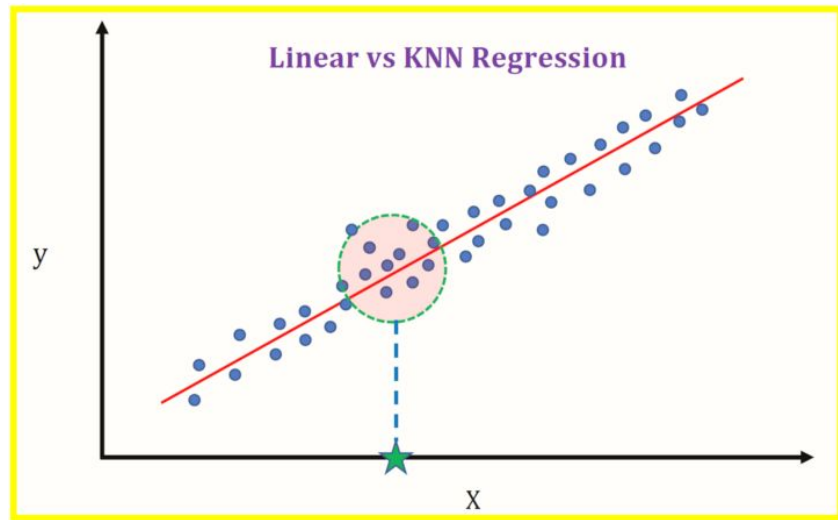
- Intuitive - each object in training has known target value
- We have  $k$  neighbors for prediction - let's average their target value!
- We can use distancing

Pros:

- Simple, not many changes from Classifier

Cons:

- All the cons of KNN



# DECISION TREE REGRESSOR

How we can change decision tree to solve regression tasks?

# DECISION TREE REGRESSOR

How we can change decision tree to solve regression tasks?

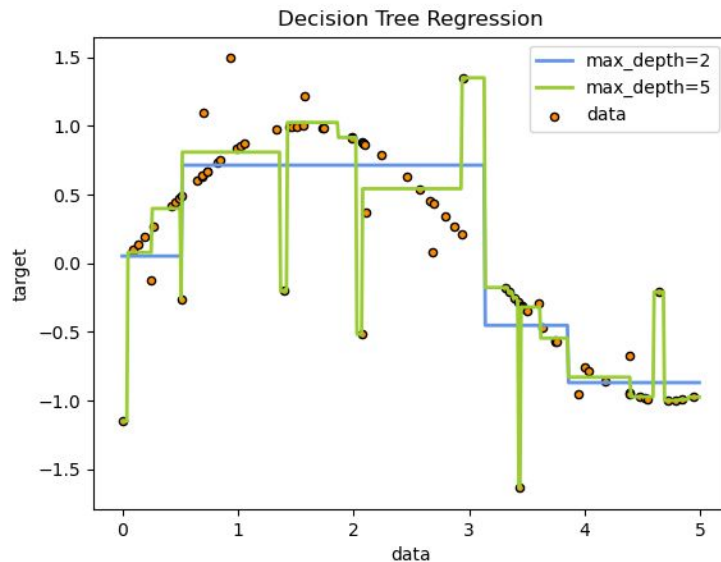
- Every leaf now contains the set of objects. Their average is the prediction we are looking for.
- We have to use other, continuous measures of information gain:
  - Variance (standard deviation)

Pros:

- Simplicity and interpretability of DT

Cons:

- Limited set of predicted values



# RANDOM FOREST REGRESSOR

How we can change random forest to solve regression task?

# RANDOM FOREST REGRESSOR

How we can change random forest to solve regression task?

- Nothing has changed, just take decision tree regressor as basic learner and average the result across estimators

Pros & cons:

- Everything is the same as in random forest classifier

# SUPPORT VECTOR MACHINE

How this is going to work?

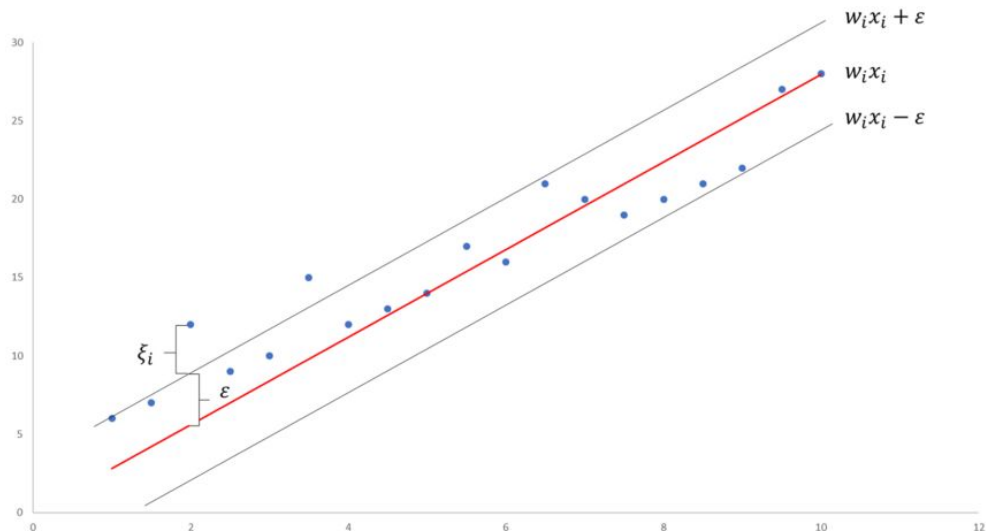
# SUPPORT VECTOR MACHINE

How this is going to work?

- Reversing the SVM task: we create the plane, as narrow as possible, which includes as many points as it can inside:

$$\text{Minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n |\xi_i|$$

$$\text{Constrain } |y_i - w_i x_i| \leq \varepsilon + |\xi_i|$$



# GRADIENT BOOSTING

How do we use GB for regression tasks?



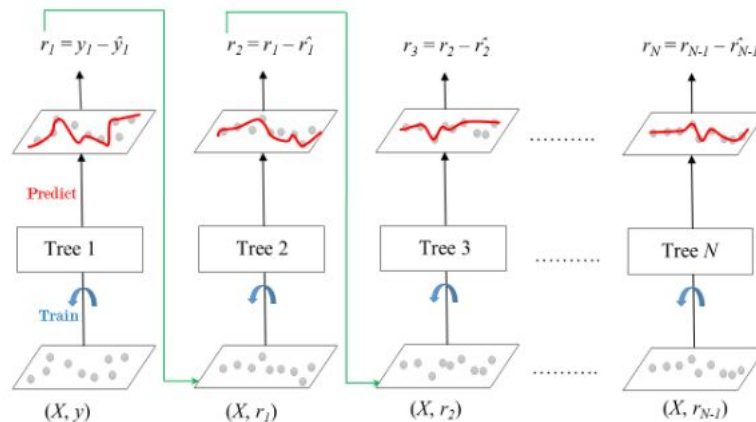
# GRADIENT BOOSTING

How do we use GB for regression tasks?

- Every new learner is fitted on error gradient with respect to ensemble of previous learners
- That means we fit every new tree on residuals from previous step

$$L_{\text{MSE}} = \frac{1}{2}(y - F(x))^2$$

$$h_m(x) = -\frac{\partial L_{\text{MSE}}}{\partial F} = y - F(x).$$



# ADVANCED HYPERPARAMETER TUNING



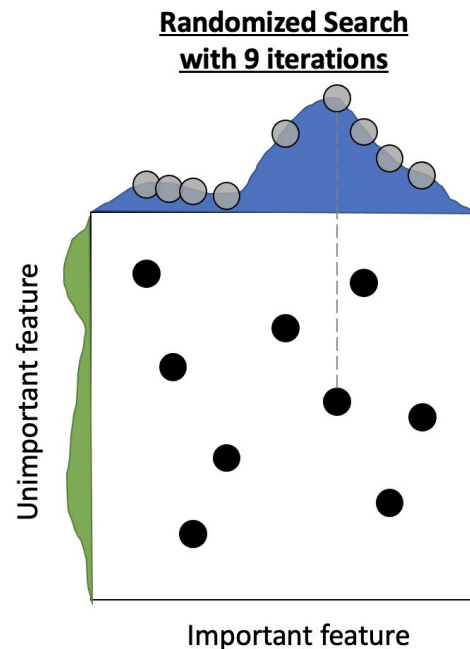
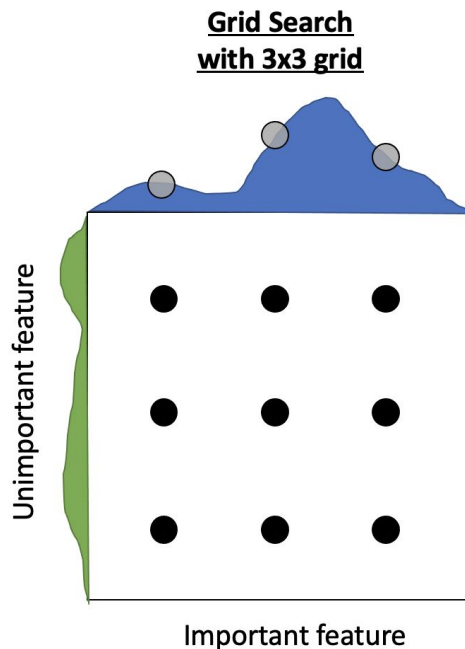
# ADVANCED HYPERPARAMETER TUNING

Which technics do you already know?

# ADVANCED HYPERPARAMETER TUNING

Which technics do you already know?

- Blind pick
- Grid Search
- Random Search



# ADVANCED HYPERPARAMETER TUNING

- HyperOpt <http://hyperopt.github.io/hyperopt/>. The idea behind can be explained through bayesian optimization

