# INTRO TO DATA SCIENCE
## LECTURE 6: REGRESSION & REGULARIZATION

LAST TIME:

- INTRO TO MACHINE LEARNING
- SUPERVISED LEARNING

QUESTIONS?

# I. REVIEW SUPERVISED LEARNING
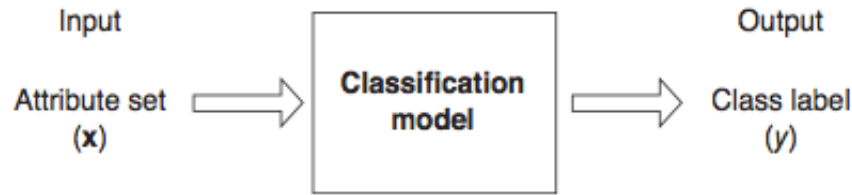# II. LINEAR REGRESSION
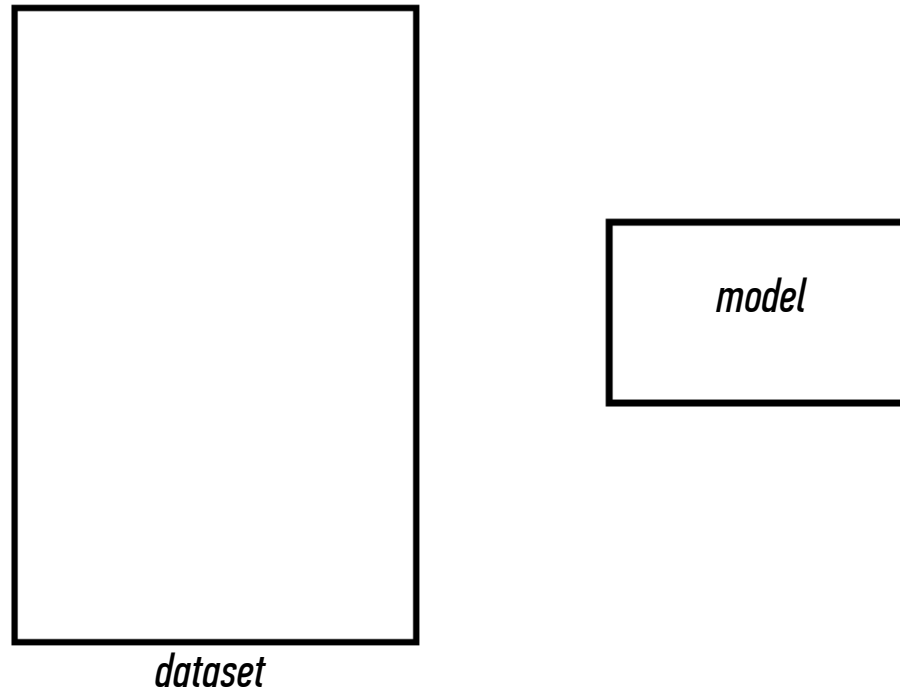# III. REGULARIZATION

# I. SUPERVISED LEARNING

*Q: How does a classification problem work?*

*A: Data in, predicted labels out.*

Input | Output

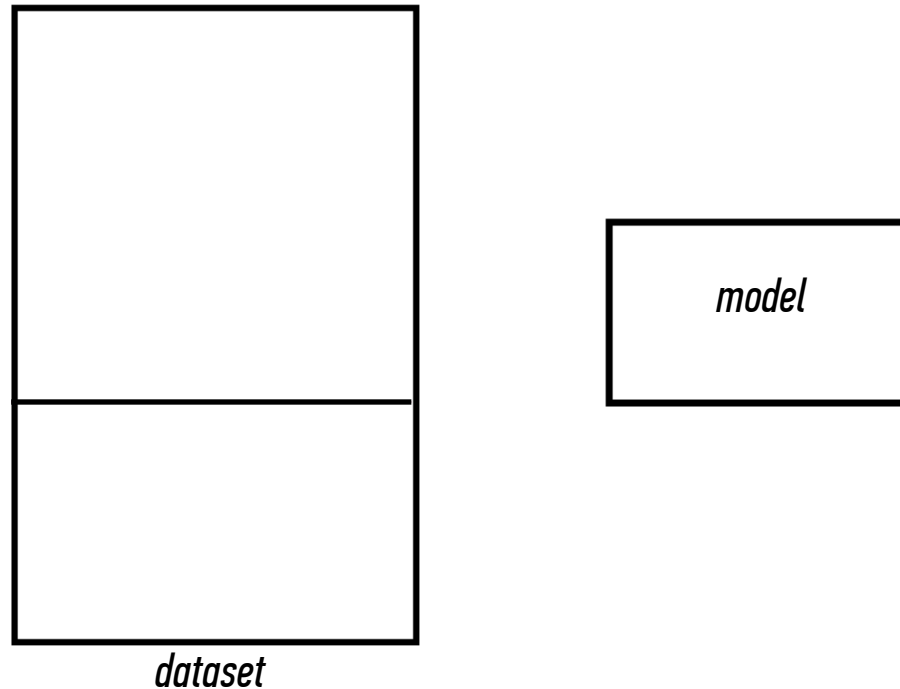Attribute set (**x**) ⇒ **Classification model** ⇒ Class label ($y$)

**Figure 4.2.** Classification as the task of mapping an input attribute set **x** into its class label $y$.
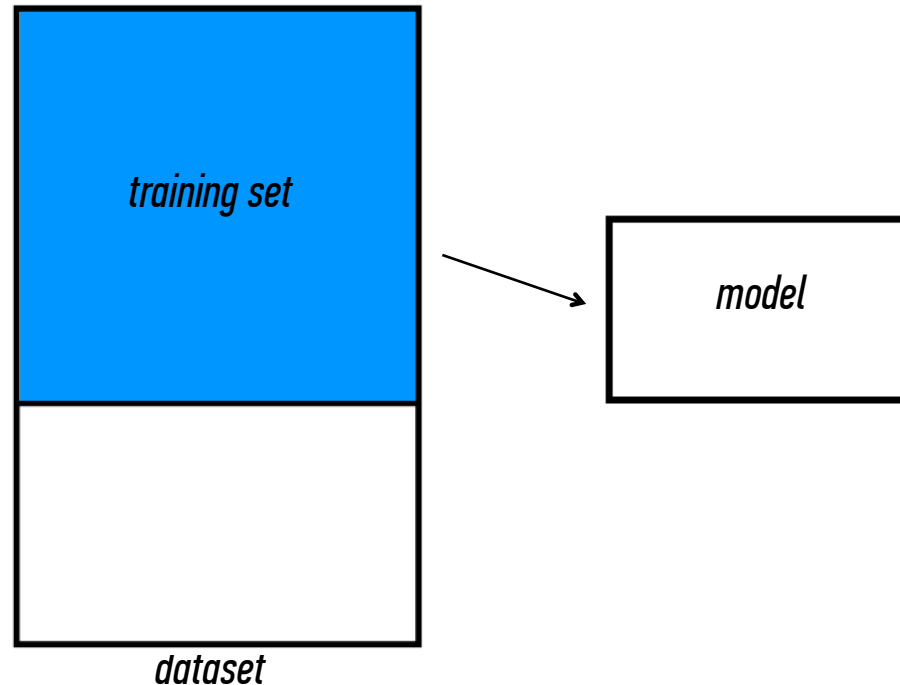
*Q: What steps does a classification problem require?*

dataset

model

*Q: What steps does a classification problem require?*

1) *split dataset*



*dataset*
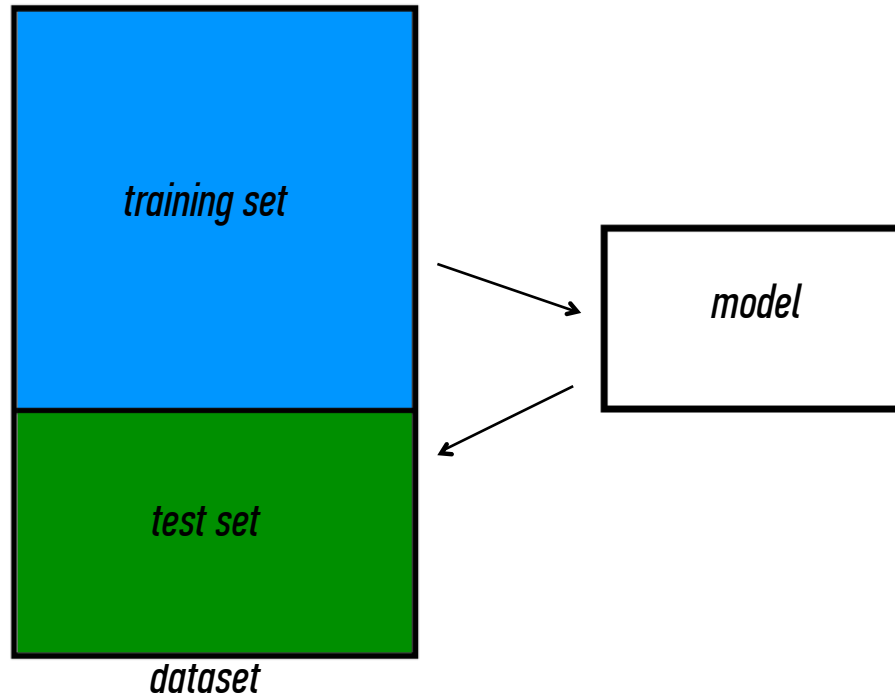
*model*

*Q: What steps does a classification problem require?*

1) split dataset

2) train model

*Q: What steps does a classification problem require?*

1) split dataset
2) train model
3) test model

*Q: What steps does a classification problem require?*

1) split dataset
2) train model
3) test model
4) make predictions

*Q: What steps does a classification problem require?*

1) split dataset
2) train model
3) test model
4) make predictions



**NOTE**

This new data is called out of sample data.

We don't know the labels for these OOS records!

training set

test set

*dataset*

model

new data

source: http://www.dtreg.com

# II. LINEAR REGRESSION

|  | *continuous* | *categorical* |
|---|---|---|
| *supervised* | ??? | ??? |
| *unsupervised* | ??? | ??? |

|  | **continuous** | **categorical** |
|---|---|---|
| **supervised** | regression | classification |
| **unsupervised** | dimension reduction | clustering |

*Q: What is a* **regression** *model?*

*Q: What is a **regression** model?*

*A: A functional relationship between input & response variables.*

*Q: What is a* **regression** *model?*

*A: A functional relationship between input & response variables*

*The* **simple linear regression** *model captures a linear relationship between a single input variable* x *and a response variable* y:

*Q: What is a **regression** model?*

*A: A functional relationship between input & response variables*

*The **simple linear regression** model captures a linear relationship between a single input variable x and a response variable y:*

$$y = \alpha + \beta x + \varepsilon$$

*Q: What do the terms in this model mean?*

$$y = \alpha + \beta x + \varepsilon$$

*Q: What do the terms in this model mean?*

$$y = \alpha + \beta x + \varepsilon$$

*A:* $y$ = **response variable** *(the one we want to predict)*

*Q: What do the terms in this model mean?*

$$y = \alpha + \beta x + \varepsilon$$

*A:*  *y =* **response variable** *(the one we want to predict)*

*x =* **input variable** *(the one we use to train the model)*

*Q: What do the terms in this model mean?*

$$y = \alpha + \beta x + \varepsilon$$

*A:*   $y =$ **response variable** *(the one we want to predict)*

    $x =$ **input variable** *(the one we use to train the model)*

    $\alpha =$ **intercept** *(where the line crosses the y-axis)*

*Q: What do the terms in this model mean?*

$$y = \alpha + \beta x + \varepsilon$$

*A:*   $y$ = **response variable** *(the one we want to predict)*

     $x$ = **input variable** *(the one we use to train the model)*

     $\alpha$ = **intercept** *(where the line crosses the y-axis)*

     $\beta$ = **regression coefficient** *(the model "parameter")*

*Q: What do the terms in this model mean?*

$$y = \alpha + \beta x + \varepsilon$$

*A:* $y$ = **response variable** *(the one we want to predict)*

$x$ = **input variable** *(the one we use to train the model)*

$\alpha$ = **intercept** *(where the line crosses the y-axis)*

$\beta$ = **regression coefficient** *(the model "parameter")*

$\varepsilon$ = **residual** *(the prediction error)*

*We can extend this model to several input variables, giving us the* **multiple linear regression** *model:*

*We can extend this model to several input variables, giving us the* **multiple linear regression** *model:*

$$y = \alpha + \beta_1 x_1 + \ldots + \beta_n x_n + \varepsilon$$

*Linear regression involves several technical assumptions and is often presented with lots of mathematical formality.*

*The math is not very important for our purposes, but you should check it out if you get serious about solving regression problems.*

*Q: How do we fit a regression model to a dataset?*

*Q: How do we fit a regression model to a dataset?*

*A: In theory, minimize the sum of the squared residuals (OLS).*

*Q: How do we fit a regression model to a dataset?*

*A: In theory, minimize the sum of the squared residuals (OLS).*

*In practice, any respectable piece of software will do this for you.*

*Q: How do we fit a regression model to a dataset?*

*A: In theory, minimize the sum of the squared residuals (OLS).*

*In practice, any respectable piece of software will do this for you.*

*But again, if you get serious about regression, you should learn how this works!*

# V. POLYNOMIAL REGRESSION

*Consider the following* **polynomial regression** *model:*

$$y = \alpha + \beta_1 x + \beta_2 x^2 + \varepsilon$$

*Consider the following* **polynomial regression** *model:*

$$y = \alpha + \beta_1 x + \beta_2 x^2 + \varepsilon$$

*Q: This represents a nonlinear relationship. Is it still a linear model?*

*Consider the following* **polynomial regression** *model:*

$$y = \alpha + \beta_1 x + \beta_2 x^2 + \varepsilon$$

*Q: This represents a nonlinear relationship. Is it still a linear model?*

*A: Yes, because it's linear in the $\beta$'s!*

*Polynomial regression allows us to fit very complex curves to data.*

$$y = \alpha + \beta_1 x + \beta_2 x^2 + \ldots + \beta_n x^n + \varepsilon$$

*Polynomial regression allows us to fit very complex curves to data.*

$$y = \alpha + \beta_1 x + \beta_2 x^2 + \ldots + \beta_n x^n + \varepsilon$$

*But there is one problem with the model we've written down so far.*

*Polynomial regression allows us to fit very complex curves to data.*

$$y \ = \ \alpha \ + \ \beta_1 x \ + \ \beta_2 x^2 \ + \ \ldots \ + \ \beta_n x^n \ + \ \varepsilon$$

*But there is one problem with the model we've written down so far.*

*Q: Does anyone know what it is?*

*Polynomial regression allows us to fit very complex curves to data.*

$$y \ = \ \alpha \ + \ \beta_1 x \ + \ \beta_2 x^2 \ + \ \ldots \ + \ \beta_n x^n \ + \ \varepsilon$$

*But there is one problem with the model we've written down so far.*

*Q: Does anyone know what it is?*

*A: This model violates one of the assumptions of linear regression!*

*This model displays* **multicollinearity,** *which means the predictor variables are highly correlated with each other.*

$$y = \alpha + \beta_1 x + \beta_2 x^2 + \dots + \beta_n x^n + \varepsilon$$

```
> x <- seq(1, 10, 0.1)
> cor(x^9, x^10)
[1] 0.9987608
```

*This model displays* **multicollinearity,** *which means the predictor variables are highly correlated with each other.*

$$y = \alpha + \beta_1 x + \beta_2 x^2 + \ldots + \beta_n x^n + \varepsilon$$

*Multicollinearity causes the linear regression model to break down, because it can't tell the predictor variables apart.*

*Q: What can we do about this?*

Q:  What can we do about this?

A:  Replace the correlated predictors with uncorrelated predictors.

*Q:  What can we do about this?*

*A:  Replace the correlated predictors with uncorrelated predictors.*

$$y = \alpha + \beta_1 f_1(x) + \beta_2 f_2(x^2) + \ldots + \beta_n f_n(x^n) + \varepsilon$$

*So far, we've seen how polynomial regression allows us to fit complex nonlinear relationships, and even to avoid multicollinearity (by using basis functions).*

*So far, we've seen how polynomial regression allows us to fit complex nonlinear relationships, and even to avoid multicollinearity (by using basis functions).*

*Q: Can a regression model be too complex?*

# V. REGULARIZATION

*Recall our earlier discussion of* **overfitting**.

*Recall our earlier discussion of* **overfitting**.

*When we talked about this in the context of classification, we said that it was a result of matching the training set too closely.*
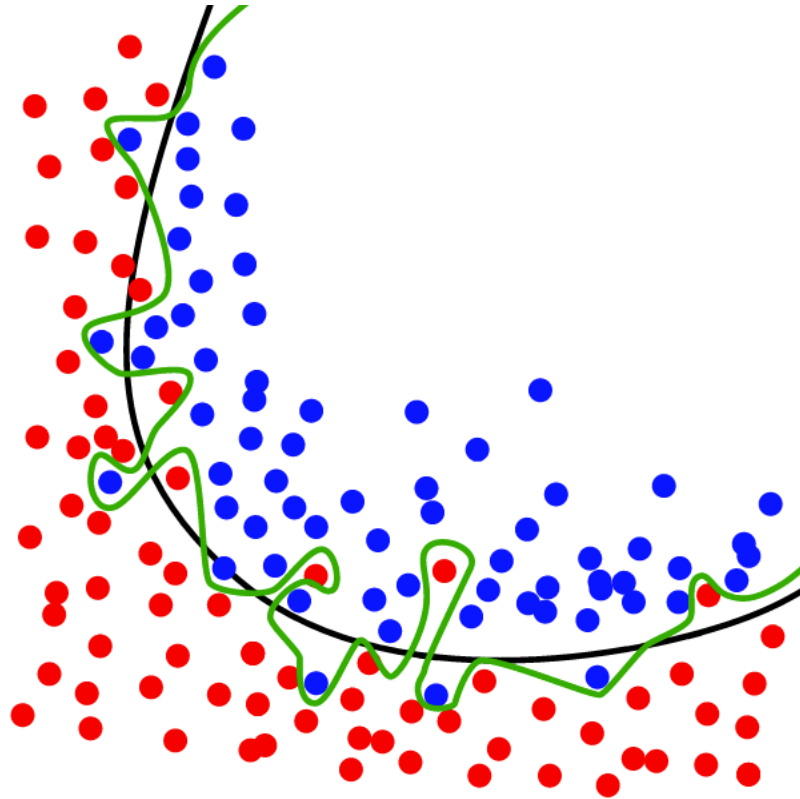
*Recall our earlier discussion of* **overfitting**.

*When we talked about this in the context of classification, we said that it was a result of matching the training set too closely.*

*In other words, an overfit model matches the* **noise** *in the dataset instead of the* **signal**.

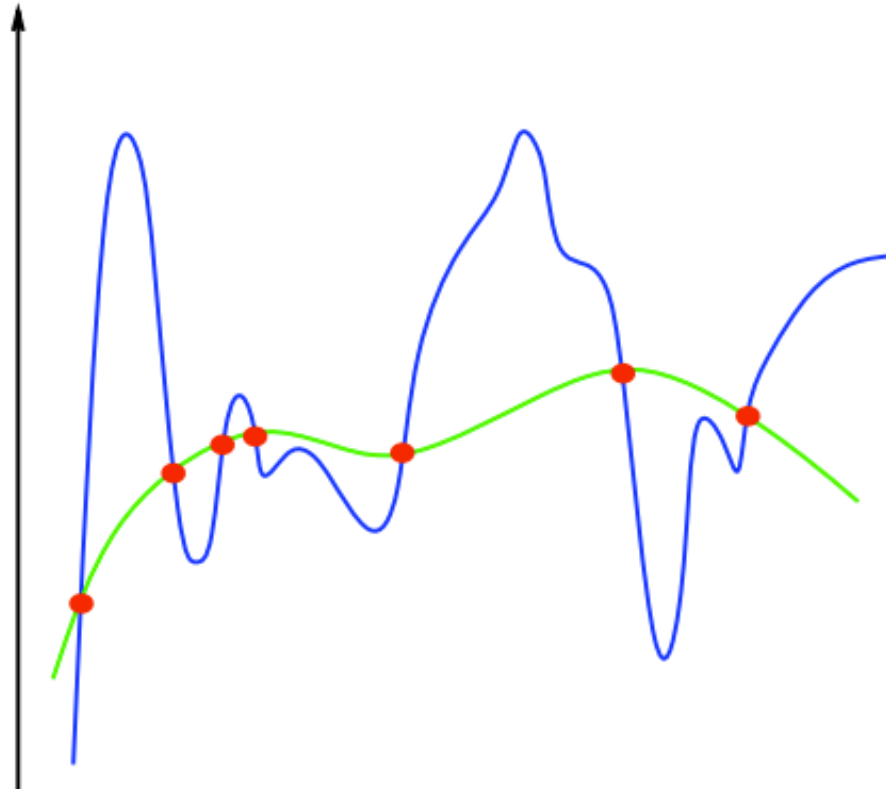source: http://upload.wikimedia.org/wikipedia/commons/1/19/Overfitting.svg

*The same thing can happen in regression.*

*It's possible to design a regression model that matches the noise in the data instead of the signal.*

*This happens when our model becomes too complex for the data to support.*

source: http://www.mit.edu/~9.520/spring12/slides/class02/class02.pdf

*Q: How do we define the **complexity** of a regression model?*

*Q: How do we define the **complexity** of a regression model?*

*A: One method is to define complexity as a function of the size of the coefficients.*

*Q:* *How do we define the* **complexity** *of a regression model?*

*A:* *One method is to define complexity as a function of the size of the coefficients.*

*Ex 1:* $\Sigma |\beta_i|$

*Ex 2:* $\Sigma \beta_i^2$

*Q: How do we define the* **complexity** *of a regression model?*

*A: One method is to define complexity as a function of the size of the coefficients.*

*Ex 1:* $\Sigma |\beta_i|$     *this is called the* **L1-norm**

*Ex 2:* $\Sigma \beta_i^2$     *this is called the* **L2-norm**

*These measures of complexity lead to the following* **regularization** *techniques:*

*These measures of complexity lead to the following* **regularization** *techniques:*

**L1 regularization**: $y = \Sigma \beta_i x_i + \varepsilon \quad st. \quad \Sigma |\beta_i| < s$

*These measures of complexity lead to the following* **regularization** *techniques:*

**L1 regularization**: $\quad y = \Sigma \beta_i x_i + \varepsilon \quad$ st. $\quad \Sigma |\beta_i| < s$

**L2 regularization**: $\quad y = \Sigma \beta_i x_i + \varepsilon \quad$ st. $\quad \Sigma \beta_i^2 < s$

*These measures of complexity lead to the following* **regularization** *techniques:*

**L1 regularization***:* $\quad y = \Sigma \, \beta_i x_i + \varepsilon \quad \text{st.} \quad \Sigma |\beta_i| < s$

**L2 regularization***:* $\quad y = \Sigma \, \beta_i x_i + \varepsilon \quad \text{st.} \quad \Sigma \, \beta_i^2 < s$

**Regularization** *refers to the method of preventing* **overfitting** *by explicitly controlling model* **complexity***.*

*These measures of complexity lead to the following* **regularization** *techniques:*

**Lasso** regularization:  $y = \Sigma \beta_i x_i + \varepsilon$   st.  $\Sigma |\beta_i| < s$

**Ridge** regularization:  $y = \Sigma \beta_i x_i + \varepsilon$   st.  $\Sigma \beta_i^2 < s$

**Regularization** *refers to the method of preventing* **overfitting** *by explicitly controlling model* **complexity**.

*Q: What problems have we seen?*

*A:*

      *1) Correlated predictor variables*

      *2) Large number of parameters allow us to overfit*

Q: What can we do about this?

A: If prediction is our only goal – nothing.

*Q:  What can we do about this?*

*A:  If prediction is our only goal – nothing.*

*Otherwise,*

> *1) Drop correlated predictors*
>
> *2) Get more data*