

# Use Machine Learning to Find Your Next Job

Samuel Taylor



## Job recommendations for 2017-09-03



assistant@samueltaylor.org

Sep 3



to sgt

Sr. Machine Learning / Artificial Intelligence Engineer @ ClosedLoop.ai - <http://www.indeed.com/cmp/ClosedLoop/jobs/Senior-Machine-Learning-f3f3a19d0d75b818>

Data Engineer @ Austin Fraser - [https://www.austinfraser.com/en-us/job/bbbh8350-data-engineer-1503529772/?utm\\_source=Indeed&utm\\_medium=organic&utm\\_campaign=Indeed](https://www.austinfraser.com/en-us/job/bbbh8350-data-engineer-1503529772/?utm_source=Indeed&utm_medium=organic&utm_campaign=Indeed)

AppSumo - Python developer @ AppSumo - [https://boards.greenhouse.io/appsumocareers/jobs/738433?gh\\_src=dognew1](https://boards.greenhouse.io/appsumocareers/jobs/738433?gh_src=dognew1)

Back-End Developer (Python) @ Beyond - [https://boards.greenhouse.io/beyond/jobs/814873?gh\\_src=ebmk7v1](https://boards.greenhouse.io/beyond/jobs/814873?gh_src=ebmk7v1)

Senior Back-End Developer @ Beyond - [https://boards.greenhouse.io/beyond/jobs/814896?gh\\_src=1xoahl1](https://boards.greenhouse.io/beyond/jobs/814896?gh_src=1xoahl1)

Software Development Principal Engineer - Austin, TX @ Dell - <https://dell.taleo.net/careersection/2/jobdetail.ftl?job=17000FQB&tz=GMT-05:00&src=JB-11346>

# Outline

- Introduction
- Asking the right question
- Gathering data
- Analysis
- Deploying

# Outline

- **Introduction**
- Asking the right question
- Gathering data
- Analysis
- Deploying

# Outline

- Introduction
- **Asking the right question**
- Gathering data
- Analysis
- Deploying

1. Have a problem

Machine learning?



# Machine learning

- Goal: Find  $f(x)$
- Problem:  $f(x)$  is unknown
- But: we can measure some points from  $f(x)$
- Algorithms to find a  $g(x)$  that approximates  $f(x)$

## Supervised

- Regression
- Classification

## Unsupervised

- Clustering

## Other stuff

- Reinforcement

## Supervised

- Regression
- Classification

## Unsupervised

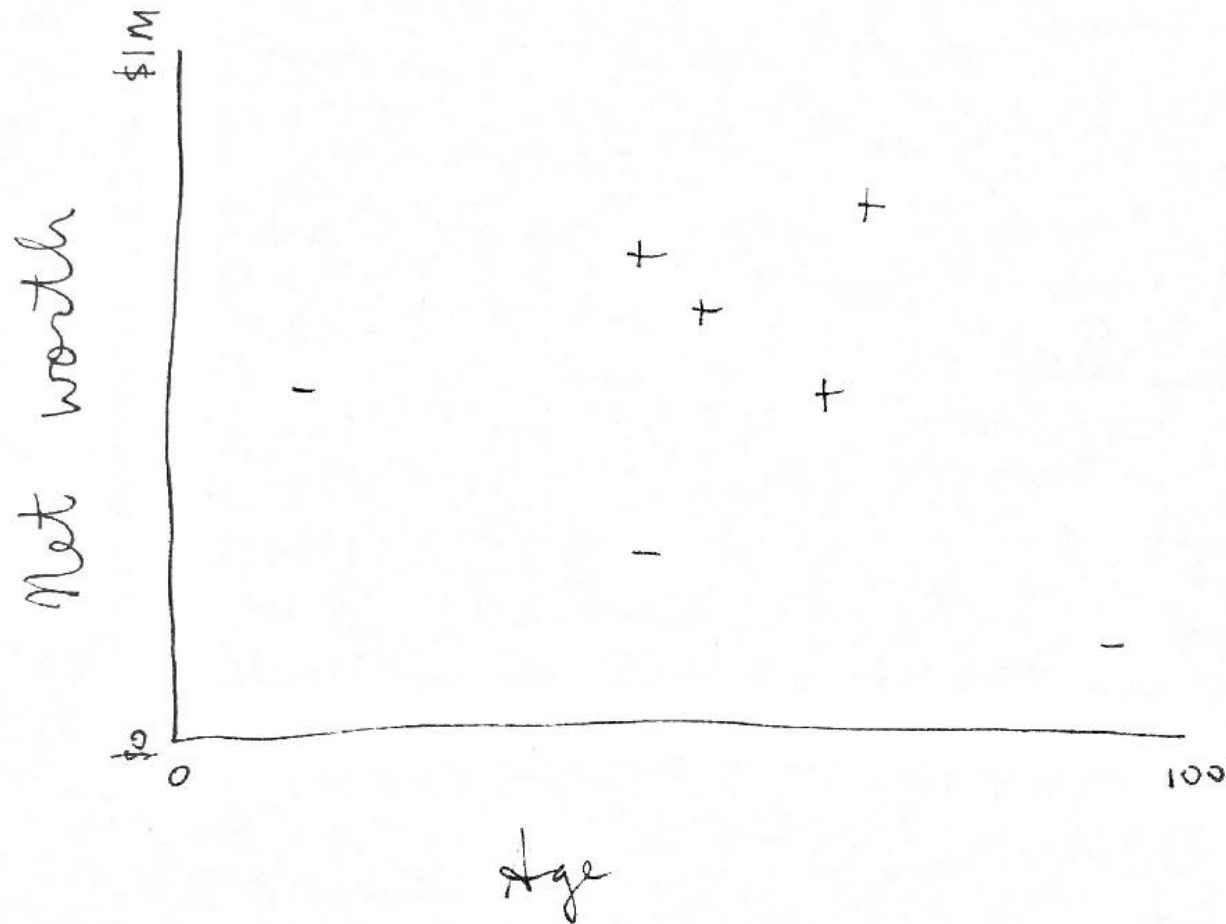
- Clustering

## Other stuff

- Reinforcement

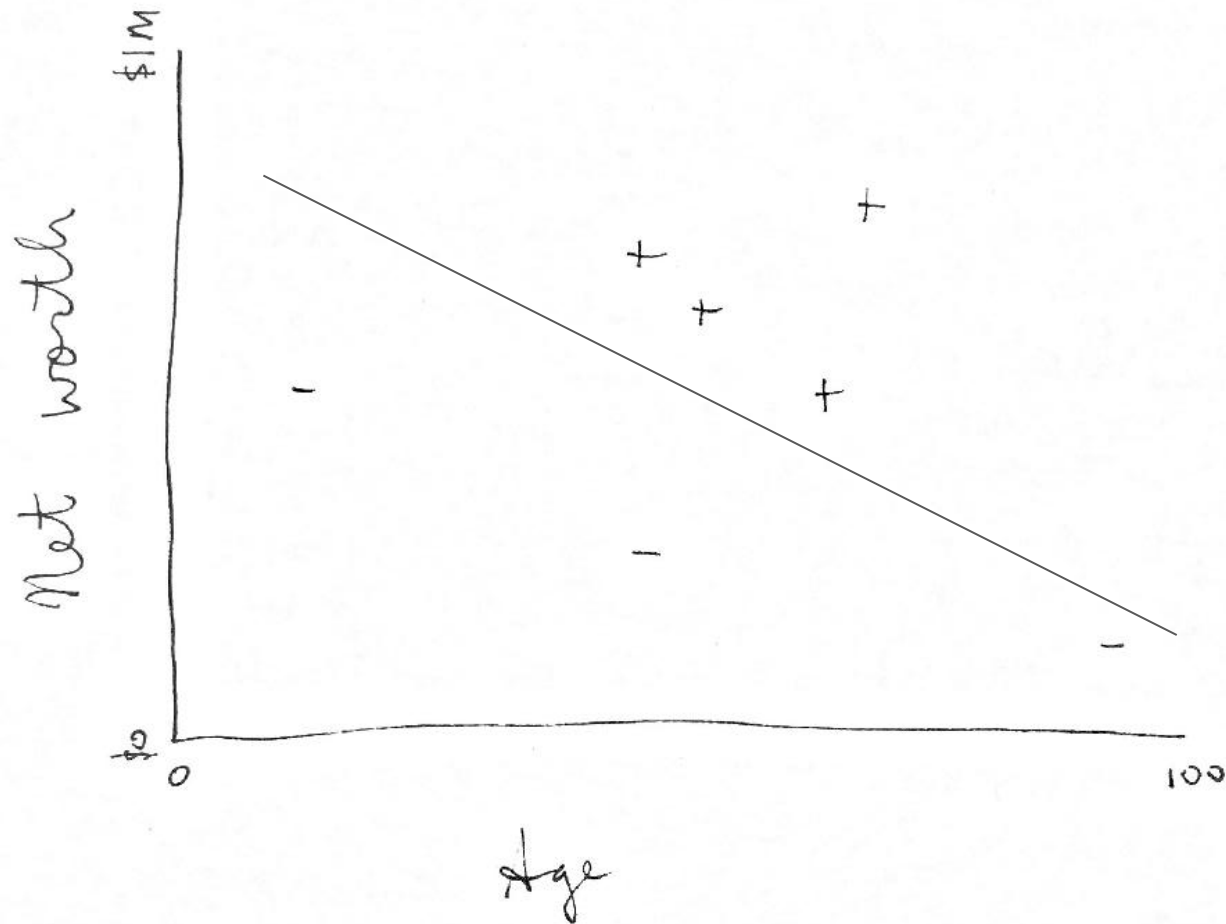
# Supervised

- Classification



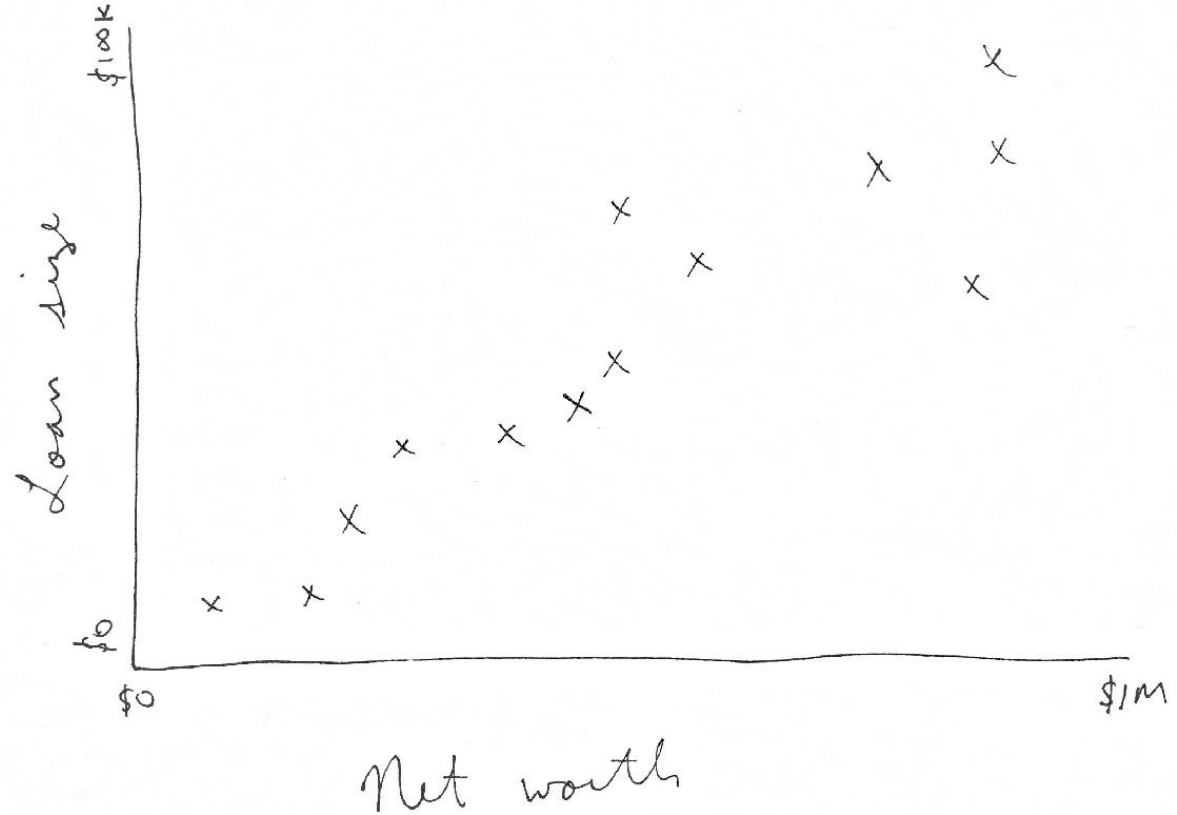
# Supervised

- Classification



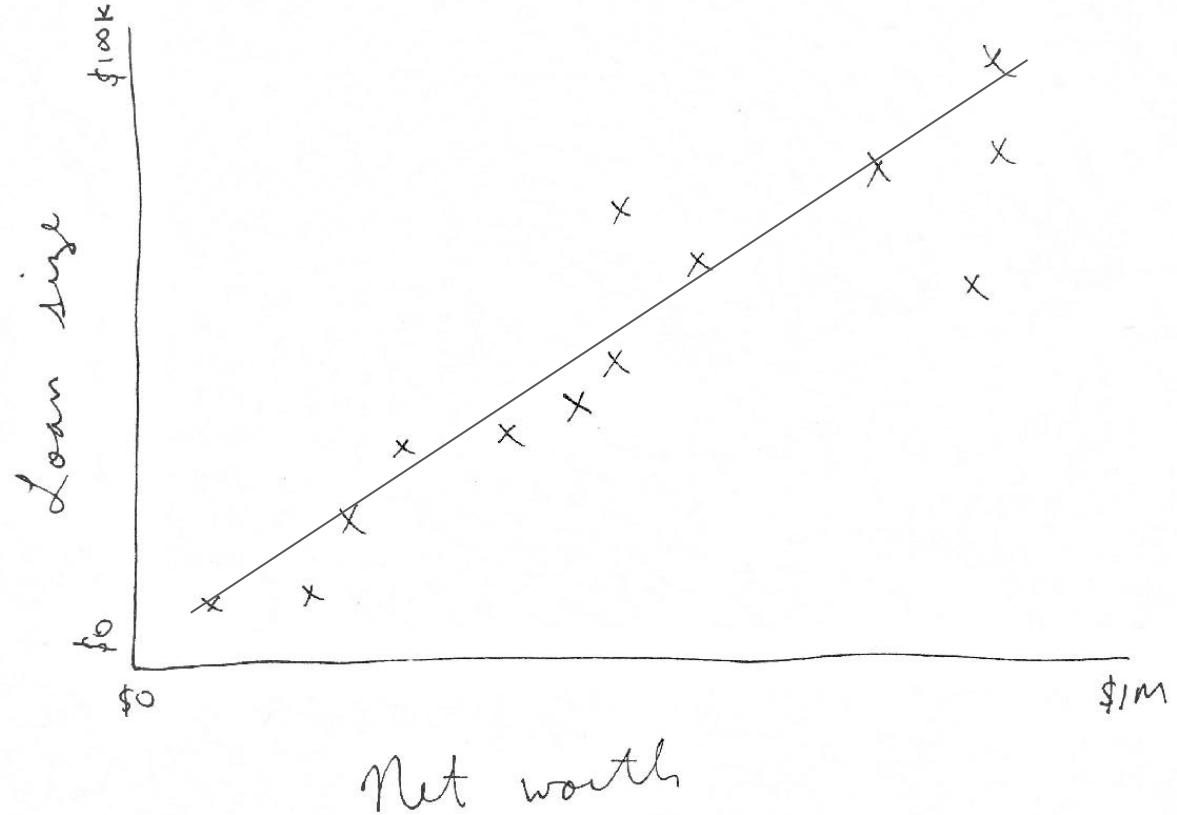
# Supervised

- Classification
- Regression



# Supervised

- Classification
- Regression



## 2. Phrase the question



# Outline

- Introduction
- Asking the right question
- **Gathering data**
- Analysis
- Deploying

# Existing

- Google it
- Government
- data.world

# Existing

- Google it
- Government
- data.world

# Create it

- Spreadsheet
- IFTTT
- Web scraping

Protip: obey robots.txt

	A	B	C	D	E
1	Title	Company	U Link		Sounds cool
2	Principal Software Architect - Austin	General Electric	/r <a href="#">Link</a>		1
3	ASIC Power Estimation Developer (Excel)	Encore Semi	/r <a href="#">Link</a>		0
4	Memory Subsystem Verification Engineer	Encore Semi	/r <a href="#">Link</a>		0
5	Senior DevOps Engineer	KIBO Software	/r <a href="#">Link</a>		0
6	Senior Manager of Software Engineering	MaxPoint	/r <a href="#">Link</a>		1
7	Data Analyst	Amherst	/r <a href="#">Link</a>		0
8	Senior Data Engineer	Visa	/r <a href="#">Link</a>		1
9	Product Development Engineer	Advanced Micro Devices, Inc.	/r <a href="#">Link</a>		0
10	Systems Analyst	Visa	/r <a href="#">Link</a>		0
11	Lead Architect - Big Data	Farmers Edge	/r <a href="#">Link</a>		1
12	Object Storage Software Engineer	IBM	/r <a href="#">Link</a>		0
13	Principal Site Reliability Engineer	Pearson	/r <a href="#">Link</a>		0
14	Senior Software Development Engineer - S	Amazon Corporate LLC	/r <a href="#">Link</a>		0
15	Systems Administrator I	University of Texas at Austin	/r <a href="#">Link</a>		0
16	Senior Database Administrator	Acxiom	/r <a href="#">Link</a>		0
17	IT Support Representative	Becker Wright Consultants	/c <a href="#">Link</a>		0
18	Software Development Engineer - Silicon C	Amazon Corporate LLC	/r <a href="#">Link</a>		0
19	Software Developer	IBM	/r <a href="#">Link</a>		0
20	Sr. Product Development Engineer	Advanced Micro Devices, Inc.	/r <a href="#">Link</a>		0
21	Front end developer	IBM	/r <a href="#">Link</a>		0
22	Full Stack Software Engineer	Indeed	/r <a href="#">Link</a>		1

	A	B	C	D	E
1	Title	Company	Link		Sounds cool
2	Principal Software Architect - Austin	General Electric	<a href="#">/r Link</a>		1
3	ASIC Power Estimation Developer (Excel-H	Encore Semi	<a href="#">/r Link</a>		0
4	Memory Subsystem Verification Engineer	Encore Semi	<a href="#">/r Link</a>		0
5	Senior DevOps Engineer	KIBO Software	<a href="#">/r Link</a>		0
6	Senior Manager of Software Engineering	MaxPoint	<a href="#">/r Link</a>		1
7	Data Analyst	Amherst	<a href="#">/r Link</a>		0
8	Senior Data Engineer	Visa	<a href="#">/r Link</a>		1
9	Product Development Engineer	Advanced Micro Devices, Inc.	<a href="#">/r Link</a>		0
10	Systems Analyst	Visa	<a href="#">/r Link</a>		0
11	Lead Architect - Big Data	Farmers Edge	<a href="#">/r Link</a>		1
12	Object Storage Software Engineer	IBM	<a href="#">/r Link</a>		0
13	Principal Site Reliability Engineer	Pearson	<a href="#">/r Link</a>		0
14	Senior Software Development Engineer - S	Amazon Corporate LLC	<a href="#">/r Link</a>		0
15	Systems Administrator I	University of Texas at Austin	<a href="#">/r Link</a>		0
16	Senior Database Administrator	Acxiom	<a href="#">/r Link</a>		0
17	IT Support Representative	Becker Wright Consultants	<a href="#">/c Link</a>		0
18	Software Development Engineer - Silicon C	Amazon Corporate LLC	<a href="#">/r Link</a>		0
19	Software Developer	IBM	<a href="#">/r Link</a>		0
20	Sr. Product Development Engineer	Advanced Micro Devices, Inc.	<a href="#">/r Link</a>		0
21	Front end developer	IBM	<a href="#">/r Link</a>		0
22	Full Stack Software Engineer	Indeed	<a href="#">/r Link</a>		1

## Existing

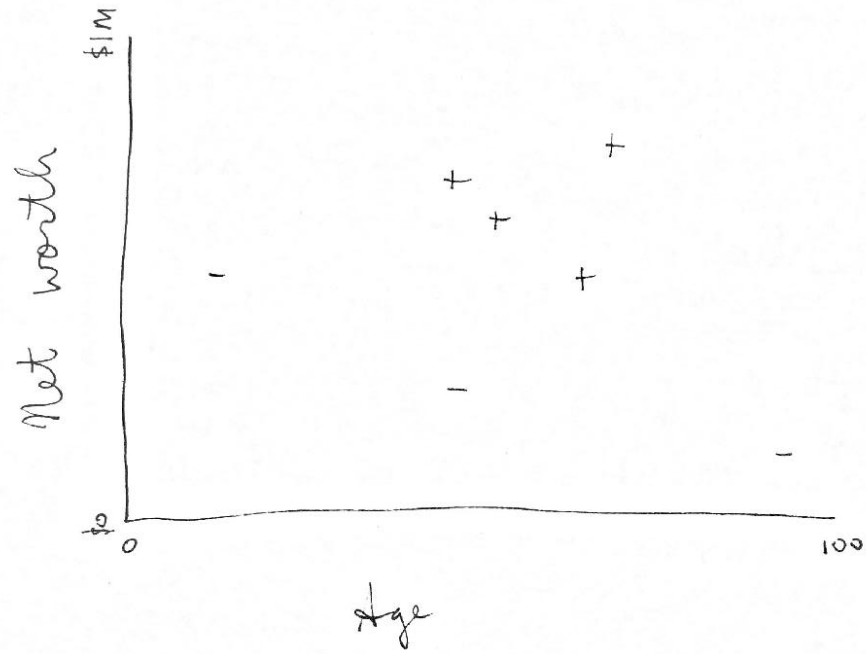
- Google it
- Government
- data.world

## Create it

- Spreadsheet
- IFTTT
- Web scraping

## Clean it

- Pandas
- scikit-learn





(Sr. Data Engineer, sounds\_cool=True)



(5, 1)

?

	Engi- neer	web	Applica- tions	sr	jr	analytics	software	data	developer
Sr. Web Applications Developer - Data Analytics	0	1	1	1	0	1	0	1	1
Jr. Software Developer	0	0	0	0	1	0	1	0	1
Sr. Data Engineer	1	0	0	1	0	0	0	1	0

(Sr. Data Engineer, sounds\_cool=True)

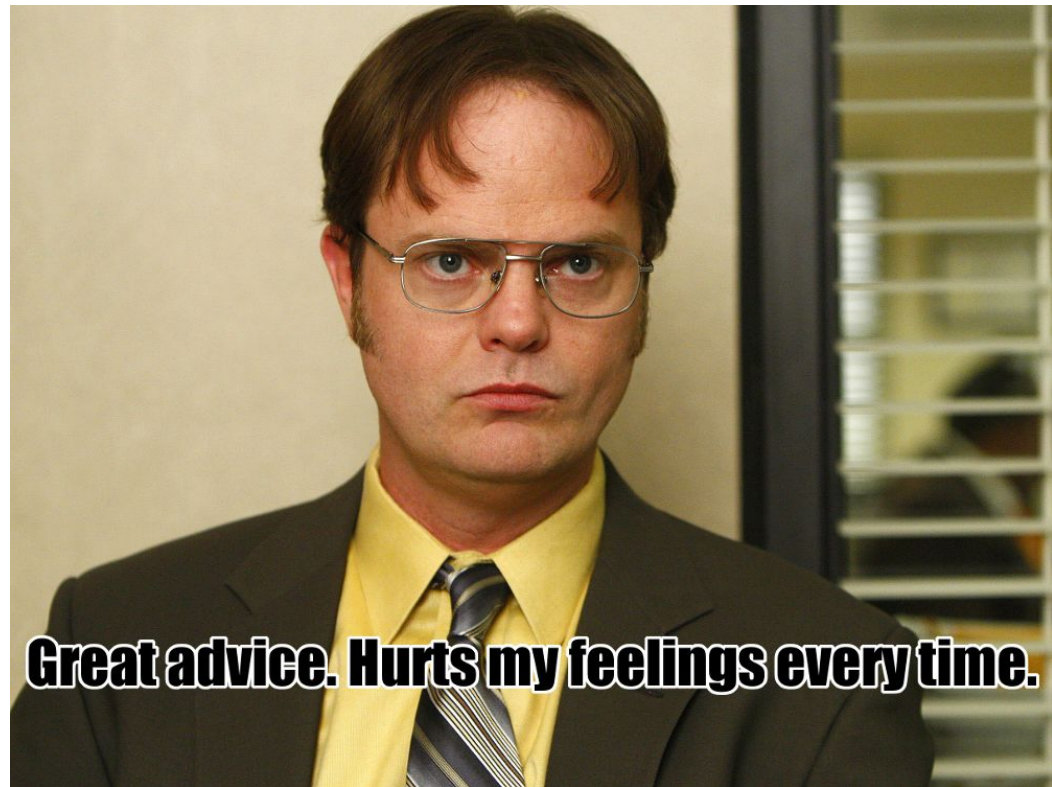


(1, 0, 0, 1, 0, 0, 0, 1, 0, 1)

# Outline

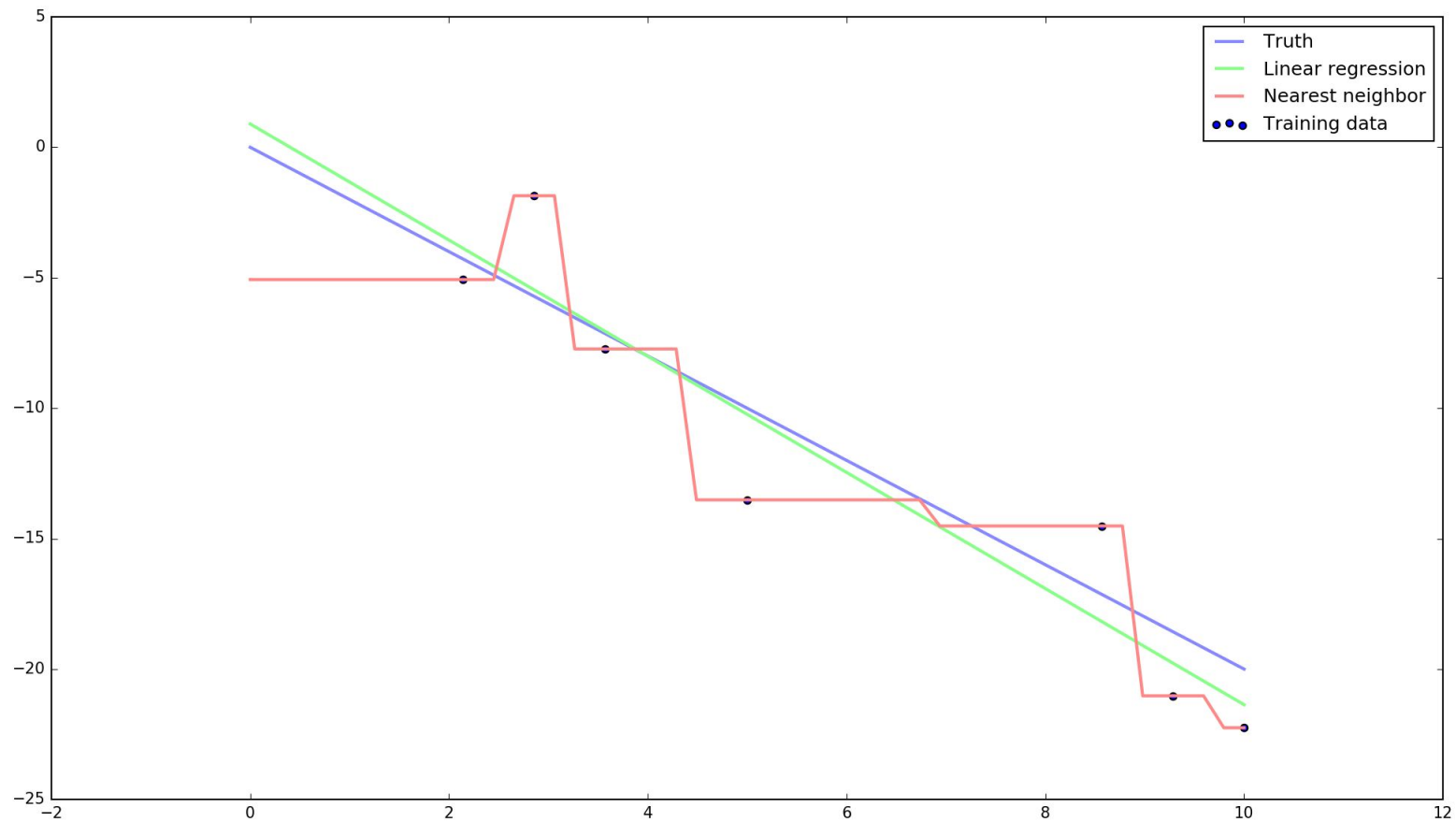
- Introduction
- Asking the right question
- Gathering data
- **Analysis**
- Deploying

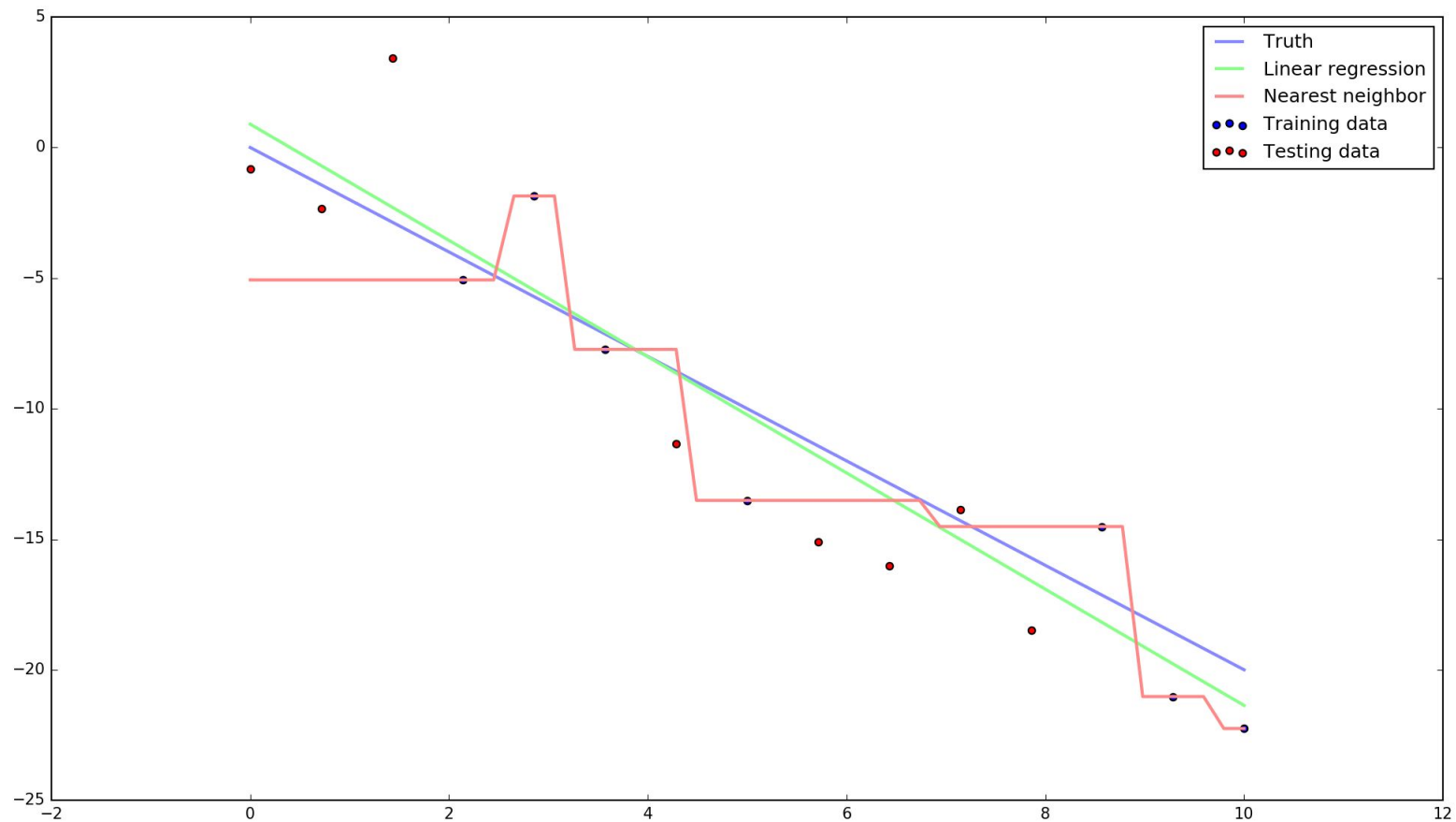
### 3. KISS



# Theory

- Approximation-generalization tradeoff







# Theory

- Approximation-generalization tradeoff
- It's just easier

# Theory

- Approximation-generalization tradeoff
- It's just easier

# Practice

- Start with simple models
  - Linear regression
  - Logistic regression

```
X = rated_jobs['title'].as_matrix()
y = rated_jobs['sounds_cool'].as_matrix()

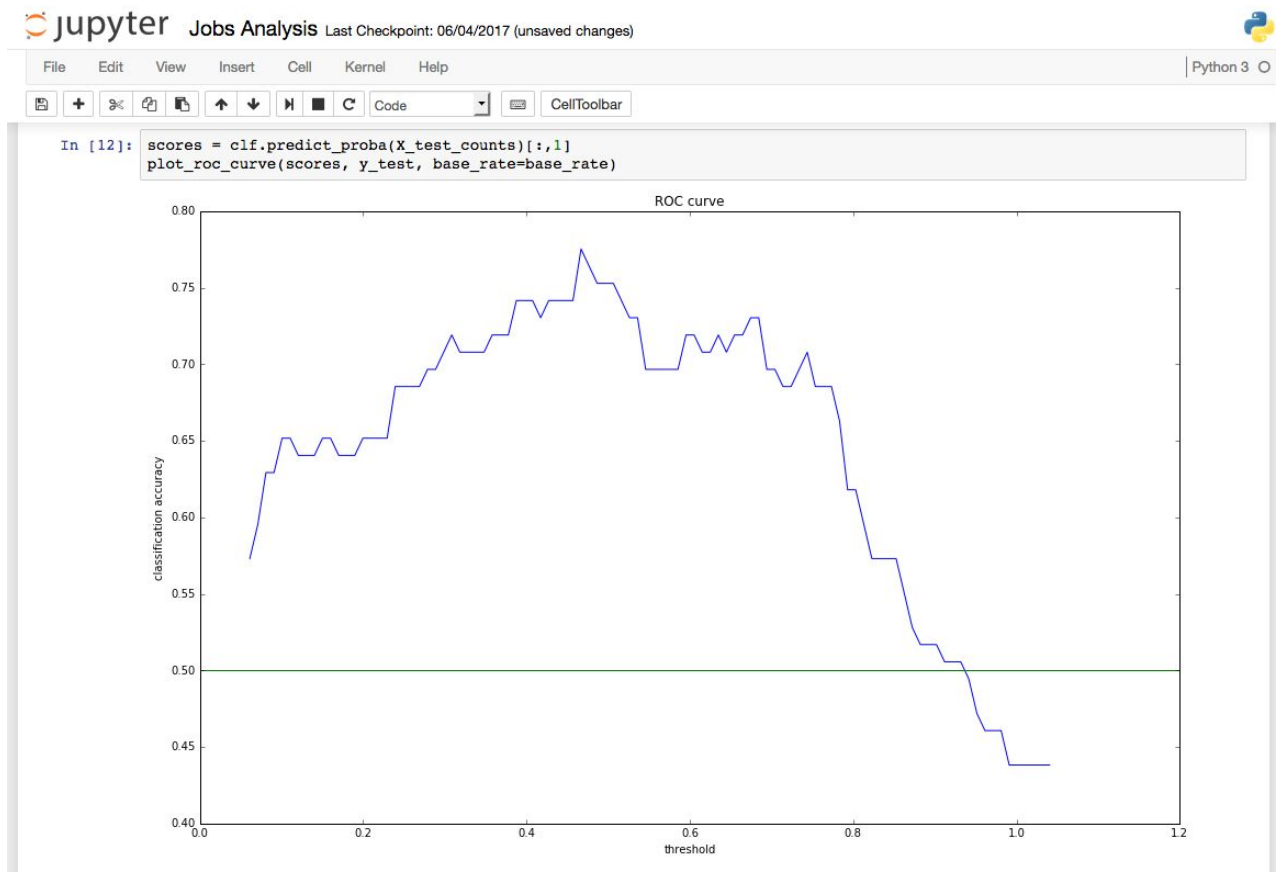
vect = CountVectorizer()
Xp = vect.fit_transform(X).toarray()
clf = LogisticRegression().fit(Xp, y)

new_job_ratings = clf.predict(new_jobs)

# array([ 0.,  0.,  0.,  1.,  0.,  0.,  0.,  1.,  0.,  0.]
```

# Recommended tools

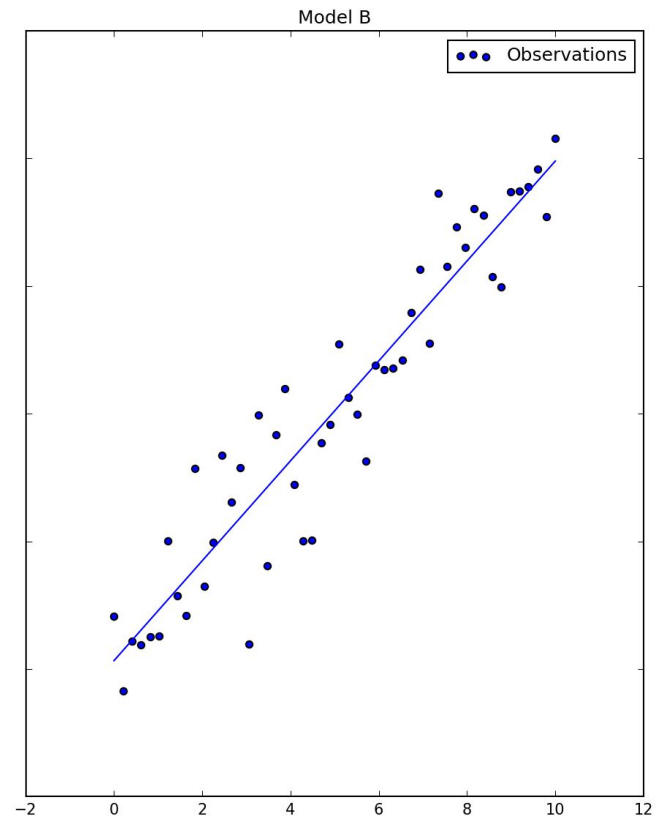
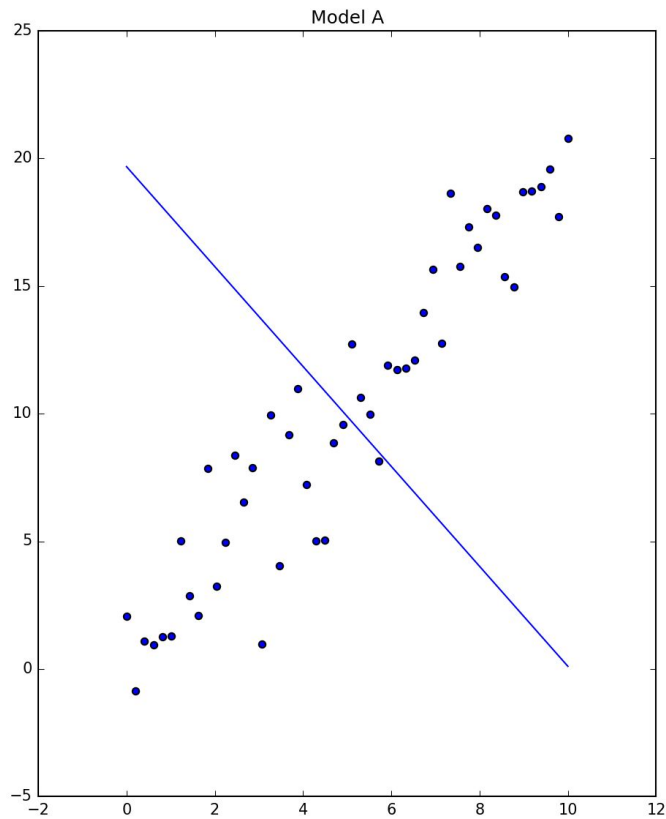
- Jupyter
- Pandas
- scikit-learn



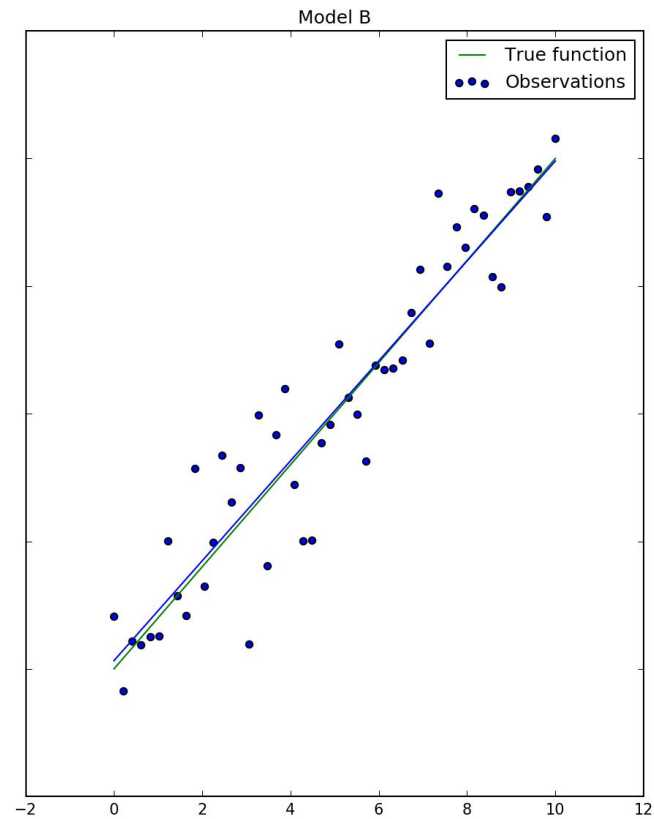
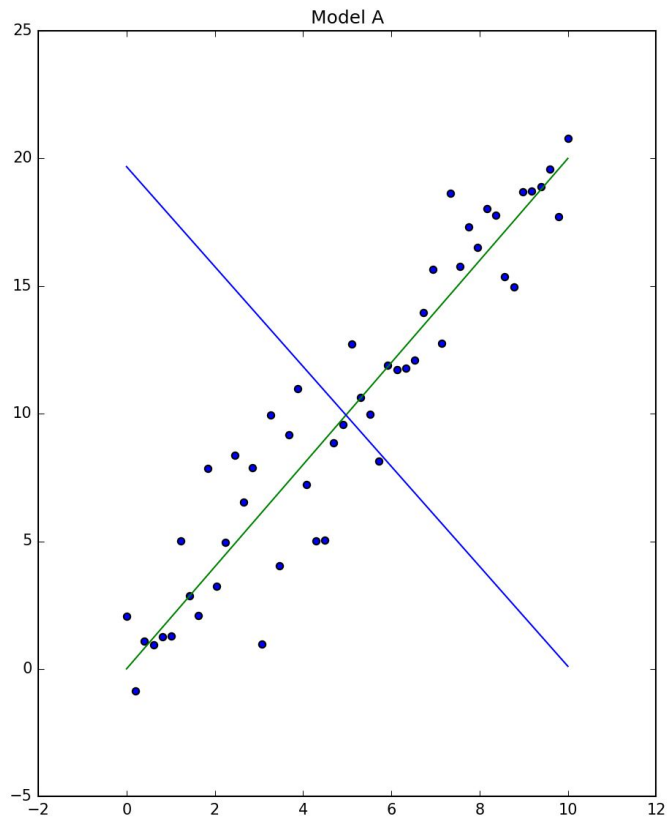
## 4. Test + iterate

How accurate is it?

# Measuring Error



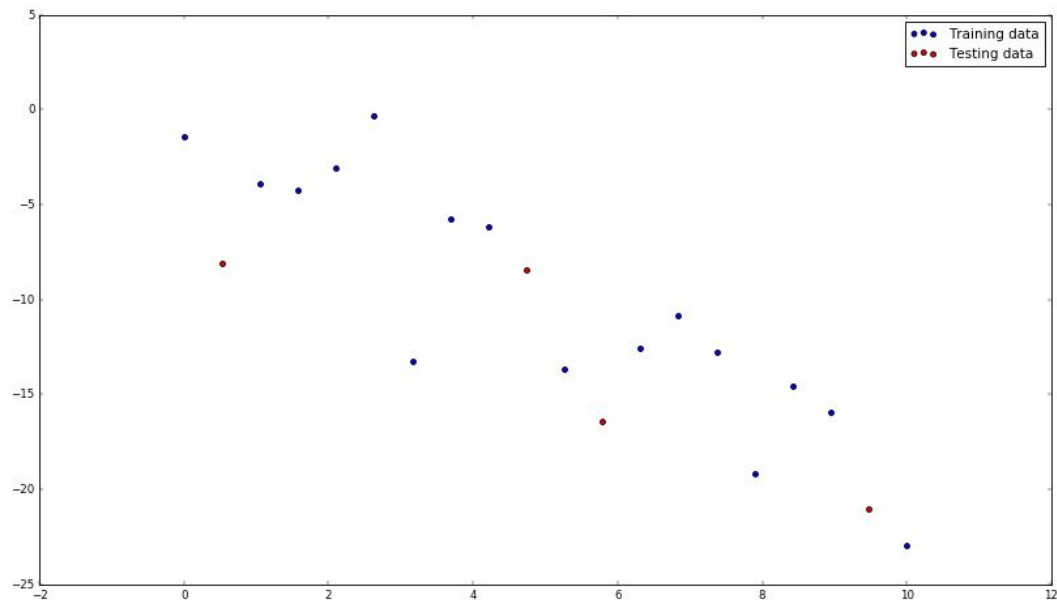
# Measuring Error





# Measuring Error

- Hold out some “testing data”

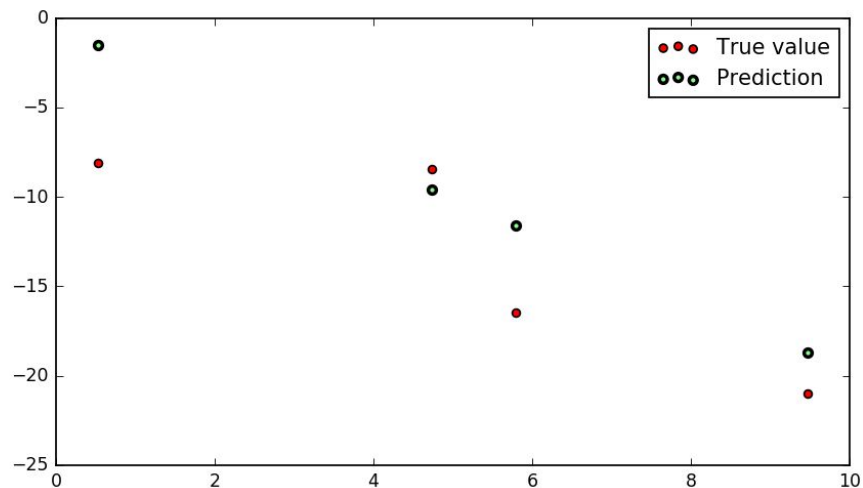


# Measuring Error

- Hold out some “testing data”
- Compare test data to prediction
- Ideally: calculate the real cost of an error
  - Cost of false positive in nuclear warhead detection: **HIGH**
  - Cost of false positive in fingerprint recognition on my phone: **SIGNIFICANTLY LOWER**

# Measuring Error

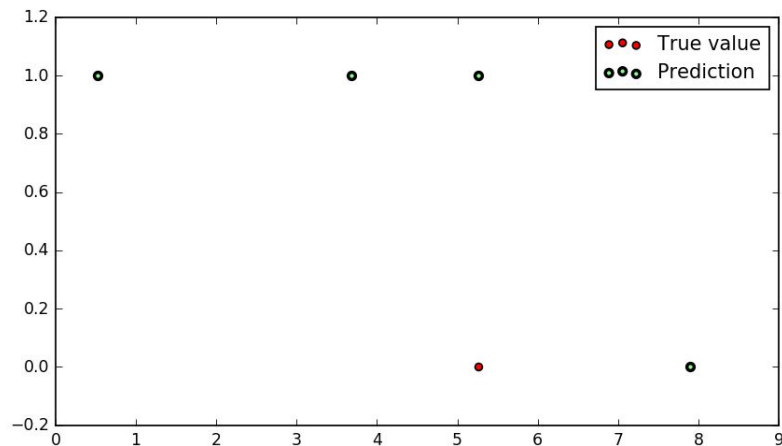
- Compare test data to prediction
- Common metric for regression: mean squared error
  - 18.35



Input	True	Predict	Diff	Sq. diff
0.53	-8.10	-1.51	-6.60	43.50
4.74	-8.47	-9.60	1.13	1.27
5.79	-16.45	-11.62	-4.83	23.30
9.47	-21.01	-18.70	-2.31	5.34

# Measuring Error

- Compare test data to prediction
- Common metric for classification: mean classification error
  - 0.25



Input	True	Predict	Error?
0.53	1	1	0
3.68	1	1	0
5.26	0	1	1
7.89	0	0	0

# Back to jobs!

- Classification error: 0.197
  - Awesome!
- But wait, it's just classifying everything as “not cool”
- Base rate for this problem is 0.197
  - No improvement

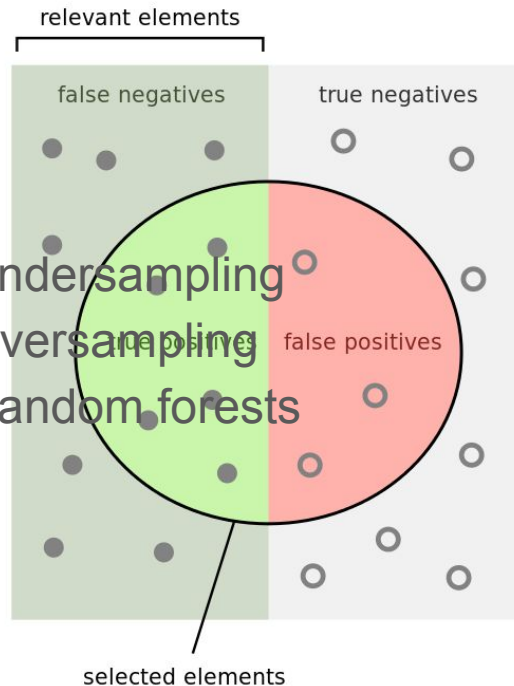
How skewed data makes me feel



# Handling imbalanced classes

- Better error metrics
  - Precision
  - Recall

- Undersampling
- Oversampling
- Random forests



How many selected items are relevant?

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

# Outline

- Introduction
- Asking the right question
- Gathering data
- Analysis
- **Deploying**



```
X = rated_jobs['title'].as_matrix()
y = rated_jobs['sounds_cool'].as_matrix()

vect = CountVectorizer()
Xp = vect.fit_transform(X).toarray()
clf = LogisticRegression().fit(Xp, y)

new_job_ratings = clf.predict(new_jobs)

# array([ 0.,  0.,  0.,  1.,  0.,  0.,  0.,  1.,  0.,  0.]
```

## Job recommendations for 2017-09-03



assistant@samueltaylor.org

Sep 3 



to sgt 

Sr. Machine Learning / Artificial Intelligence Engineer @ ClosedLoop.ai - <http://www.indeed.com/cmp/ClosedLoop/jobs/Senior-Machine-Learning-f3f3a19d0d75b818>

Data Engineer @ Austin Fraser - [https://www.austinfraser.com/en-us/job/bbbh8350-data-engineer-1503529772/?utm\\_source=Indeed&utm\\_medium=organic&utm\\_campaign=Indeed](https://www.austinfraser.com/en-us/job/bbbh8350-data-engineer-1503529772/?utm_source=Indeed&utm_medium=organic&utm_campaign=Indeed)

AppSumo - Python developer @ AppSumo - [https://boards.greenhouse.io/appsumocareers/jobs/738433?gh\\_src=dognew1](https://boards.greenhouse.io/appsumocareers/jobs/738433?gh_src=dognew1)

Back-End Developer (Python) @ Beyond - [https://boards.greenhouse.io/beyond/jobs/814873?gh\\_src=ebmk7v1](https://boards.greenhouse.io/beyond/jobs/814873?gh_src=ebmk7v1)

Senior Back-End Developer @ Beyond - [https://boards.greenhouse.io/beyond/jobs/814896?gh\\_src=1xoahl1](https://boards.greenhouse.io/beyond/jobs/814896?gh_src=1xoahl1)

Software Development Principal Engineer - Austin, TX @ Dell - <https://dell.taleo.net/careersection/2/jobdetail.ftl?job=17000FQB&tz=GMT-05:00&src=JB-11346>

## 4. Test + iterate

# The Recruiter Repellant 3000

<demo>

4. Test + iterate  
+ iterate  
+ iterate



1. Have a problem
2. Phrase the question
3. Try the simplest thing
4. Test and iterate

# More resources

- [Learning from Data](#)
- [Practical Business Python](#)



The logo for Indeed, featuring a blue stylized 'i' with a curved line above it, followed by the word 'indeed' in a blue sans-serif font, and a registered trademark symbol (®) to the right.

indeed®

# Samuel Taylor

sgt@samuelstaylor.org

@SamuelDataT

