# Profit Curves and Imbalanced Classes

Darren Reger Lecture for Galvanize DSI

# I have 5 models, which is the best?

- Model 1 with Accuracy 0.977

- Model 2 with Accuracy 0.02

- Model 3 with Accuracy 0.98

- Model 4 with Accuracy 0.88

- Model 5 with Accuracy 0.748

# What about now?

- Model 1 with Precision 0.44 and Recall 0.6

- Model 2 with Precision 0.02 and Recall 1.0

- Model 3 with Precision 0 and Recall 0

- Model 4 with Precision 0.115 and Recall 0.75

- Model 5 with Precision 0.0672 and Recall 0.9

# Does this help?

|  | Predicted: Yes | Predicted: No |
|---|---|---|
| Actual: Yes | 12 | 15 |
| Actual: No | 8 | 965 |

|  | Predicted: Yes | Predicted: No |
|---|---|---|
| Actual: Yes | 20 | 980 |
| Actual: No | 0 | 0 |

|  | Predicted: Yes | Predicted: No |
|---|---|---|
| Actual: Yes | 0 | 0 |
| Actual: No | 20 | 980 |

|  | Predicted: Yes | Predicted: No |
|---|---|---|
| Actual: Yes | 15 | 115 |
| Actual: No | 5 | 865 |

|  | Predicted: Yes | Predicted: No |
|---|---|---|
| Actual: Yes | 18 | 250 |
| Actual: No | 2 | 730 |

# Discussion
# of
# Business
# Applications

# Revisiting Confusion Matrix

$$\begin{pmatrix} TP & FP \\ FN & TN \end{pmatrix} \rightarrow \begin{pmatrix} \dfrac{TP}{TP+FN} & \dfrac{FP}{FP+TN} \\ \dfrac{FN}{TP+FN} & \dfrac{TN}{FP+TN} \end{pmatrix}$$

$$\begin{pmatrix} \dfrac{TP}{TP+FN}\,P_+ & \dfrac{FP}{FP+TN}\,P_- \\ \dfrac{FN}{TP+FN}\,P_+ & \dfrac{TN}{FP+TN}\,P_- \end{pmatrix}$$

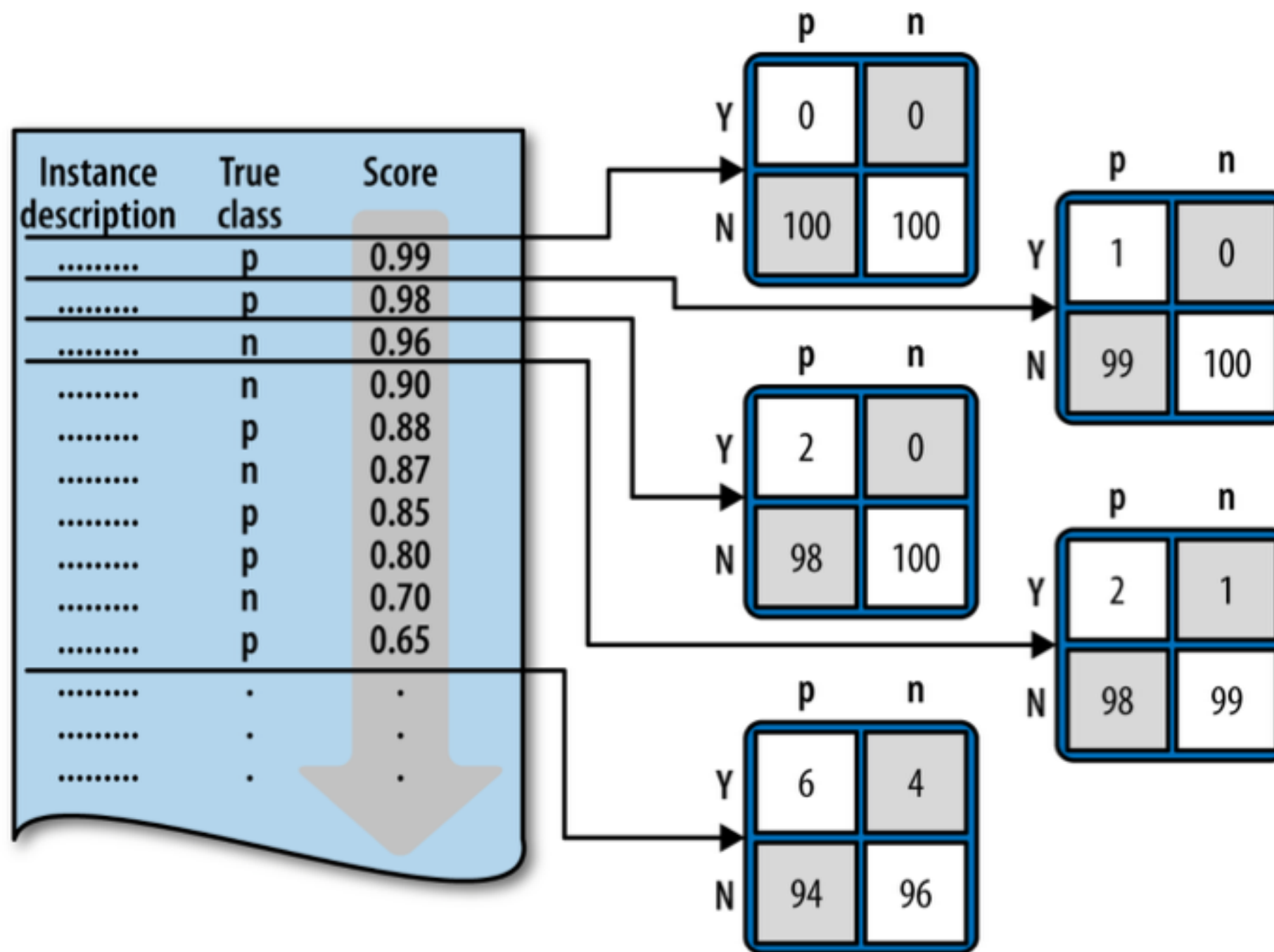$$\text{Profit Matrix} = \begin{pmatrix} B_{P_+} & C_{P_+} \\ C_{P_-} & B_{P_-} \end{pmatrix}$$

# Profit Calculation

Combining information from the **Confusion matrix** and the **Cost-Benefit matrix** we can calculate **Expected Profit!**

|  | Actual | |
|---|---|---|
| | p | n |
| **Y** | b(Y,p) | c(Y,n) |
| **N** | c(N,p) | b(N,n) |

(Predicted)

$$
\begin{aligned}
E[Profit] \; = \; & P(Y,p) \cdot b(Y,p) + P(Y,n) \cdot c(Y,n) + \\
& P(N,p) \cdot c(N,p) + P(N,n) \cdot b(N,n) \\[8pt]
= \; & P(Y|p) \cdot P(p) \cdot b(Y,p) + P(Y|n) \cdot P(n) \cdot c(Y,n) + \\
& P(N|p) \cdot P(p) \cdot c(N,p) + P(N|n) \cdot P(n) \cdot b(N,n) \\[8pt]
= \; & P(p) \cdot [P(Y|p) \cdot b(Y,p) + P(N|p) \cdot c(N,p)] + \\
& P(n) \cdot [P(Y|n) \cdot c(Y,n) + P(N|n) \cdot b(N,n)]
\end{aligned}
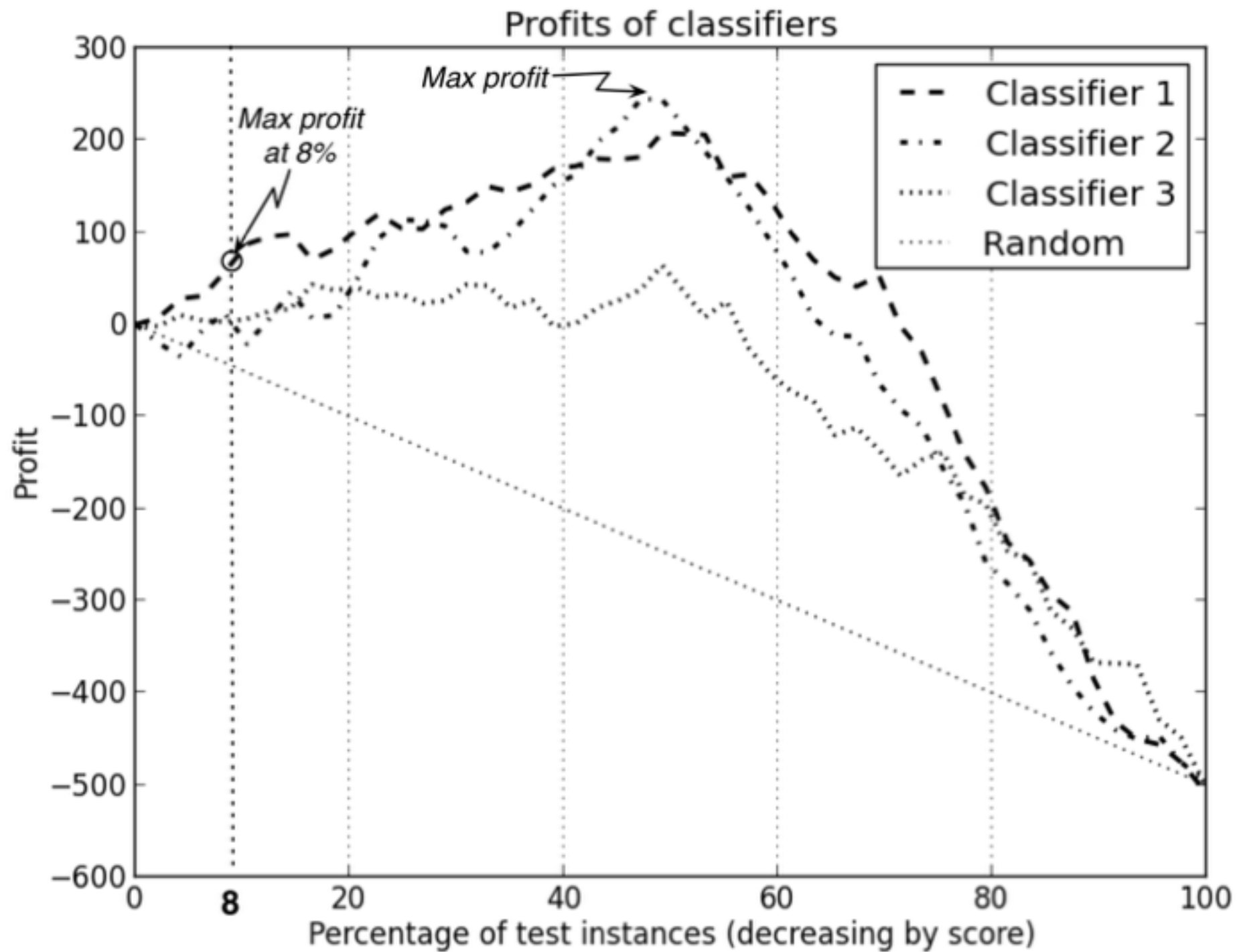$$

# Let's Make a Curve

# Making Profit Curves

For a given model $f$, each threshold value $T$ gives a point on the Profit Curve

Model score is the threshold probability classifying $+$ vs $-$

1. Allow $T$ to be the maximum score
2. $TP = 0, FP = 0$
3. Calculate $E[Profit]$
4. For each observation, $i$:
   - If $\hat{\pi}_i > T \longrightarrow$ increment TP
   - Else $\longrightarrow$ increment FP
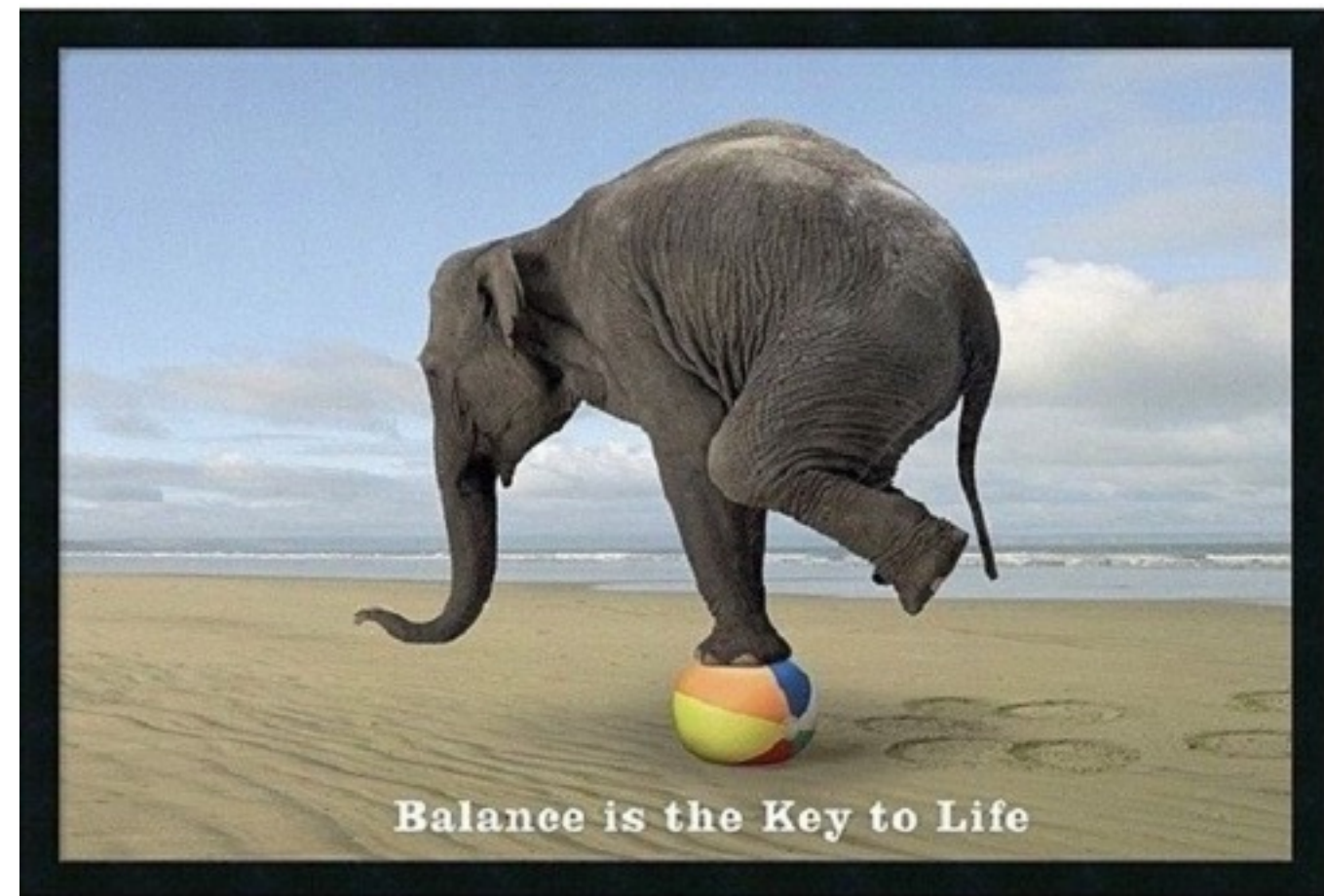5. Add point (% Test Instances predicted Positive, $E[Profit]$) to the Profit Graph

Increment $T$ from max-score to min-score, repeating steps 1-4

# $$$$$$$$$



Profits of classifiers

# Imbalanced Data

- Where do we see it?

- Why is it bad?

- What can we do?

  - Assign Weights

  - Balancing Classes



Balance is the Key to Life
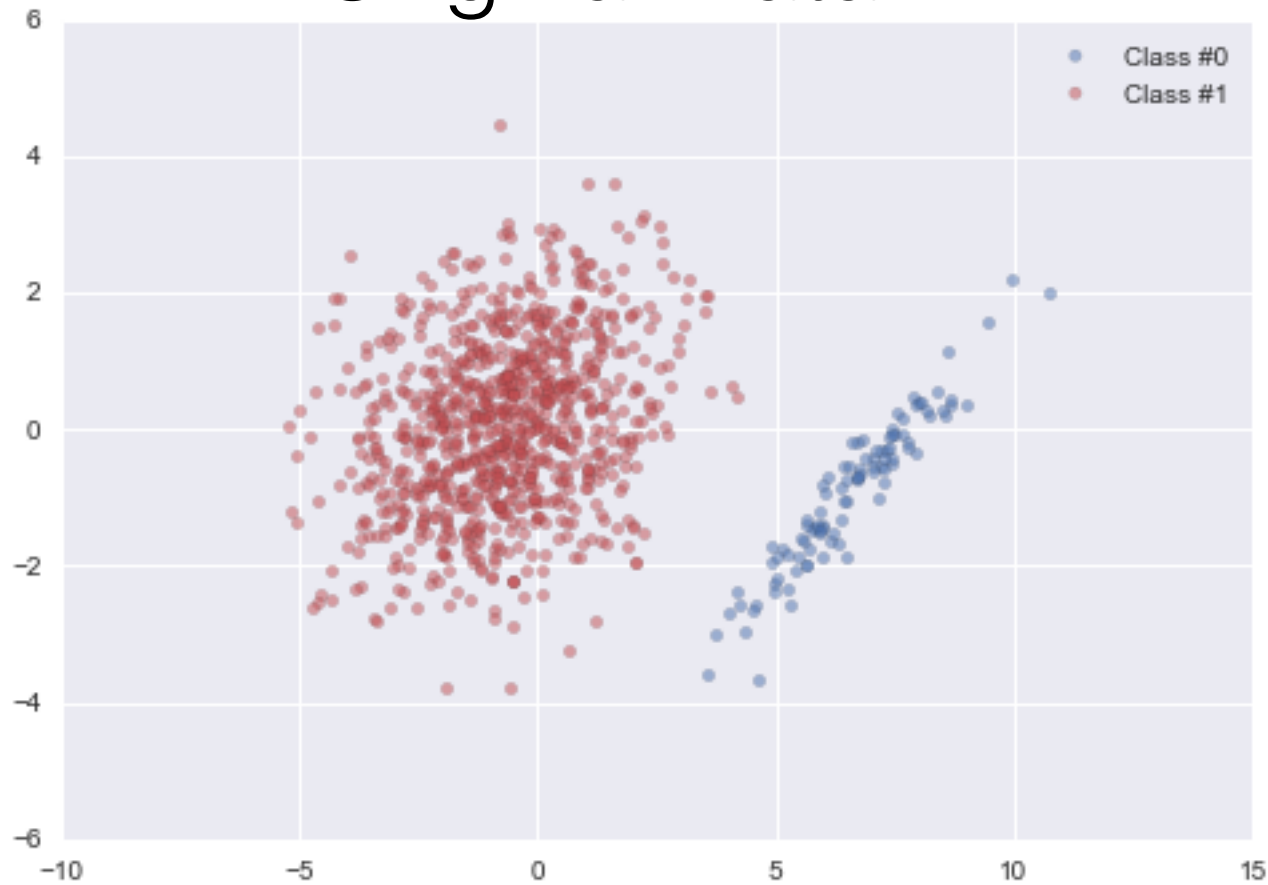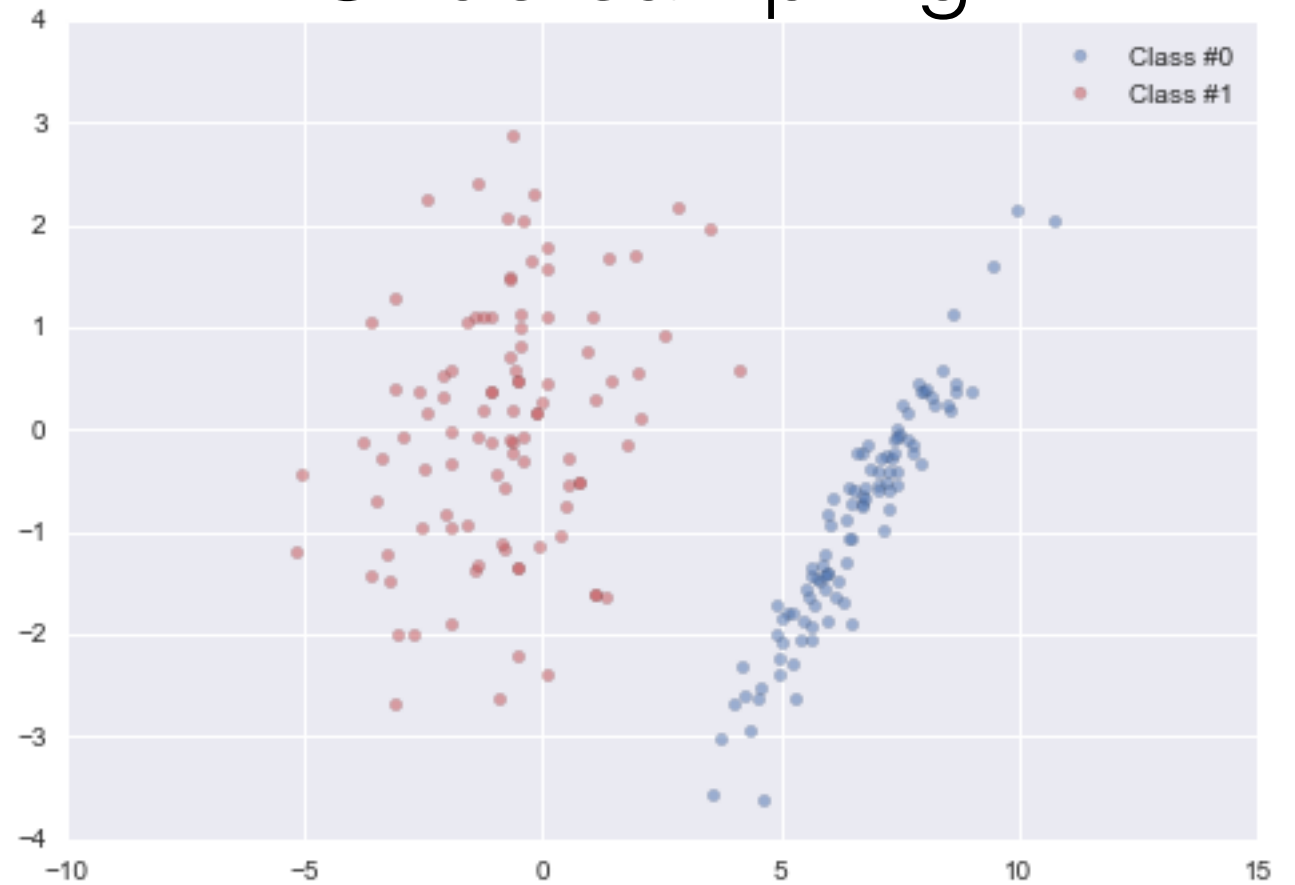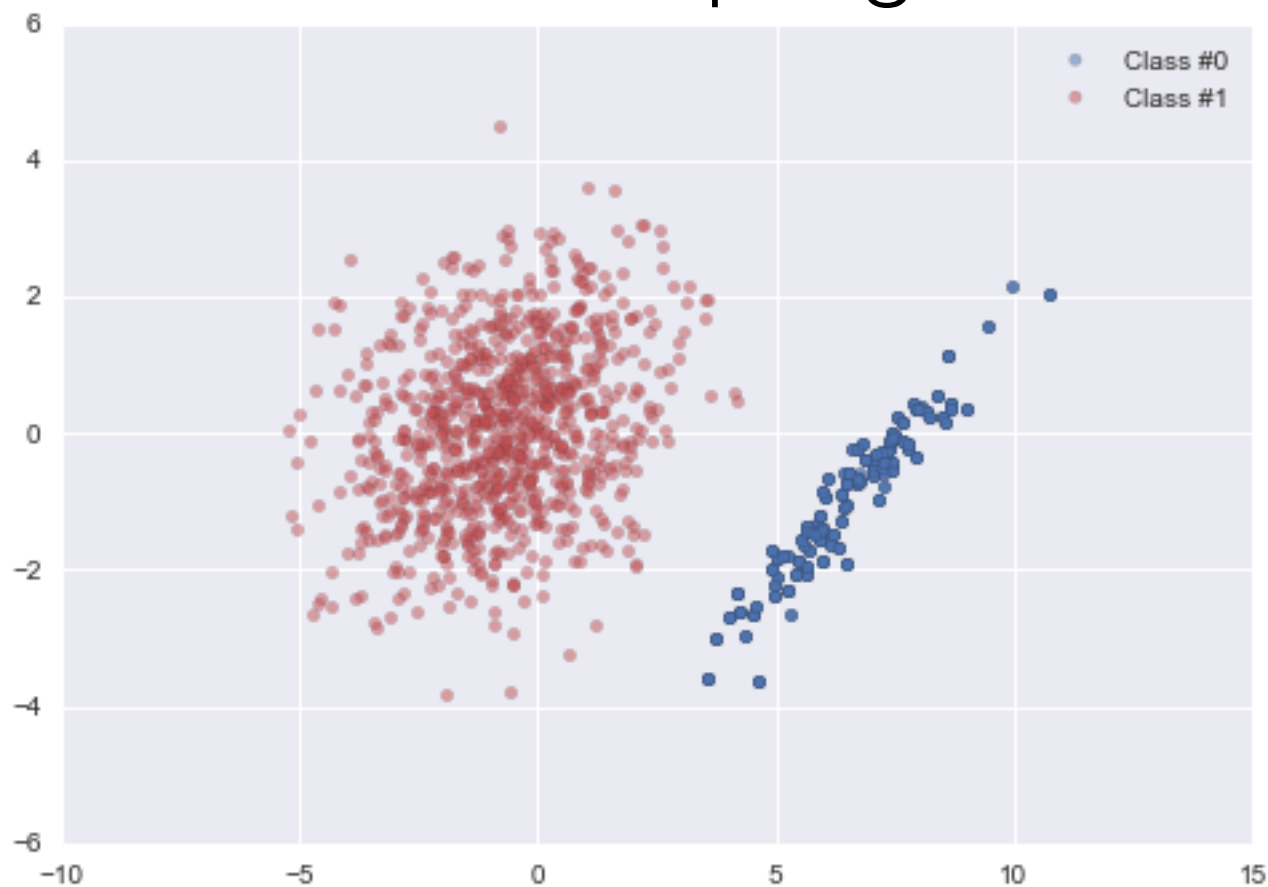
Ateeq Ahmed Siddiqui

# Choosing Cutoff Point



0.064 (Spec = 0.76, Sens = 0.64)

0.100 (Spec = 0.84, Sens = 0.44)

0.300 (Spec = 0.96, Sens = 0.15)

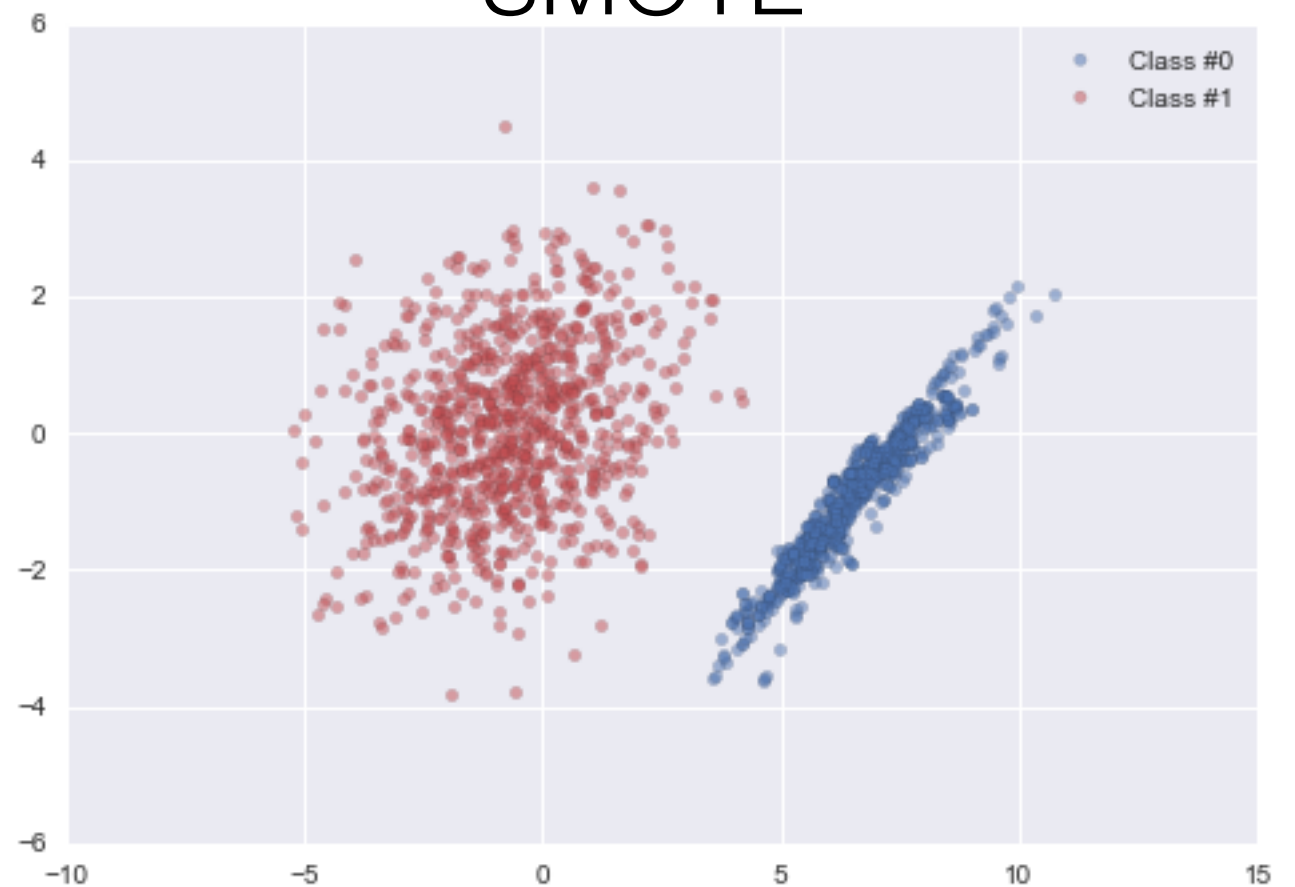0.500 (Spec = 0.99, Sens = 0.07)

**Original Data**

**Undersampling**

**Oversampling**

**SMOTE**

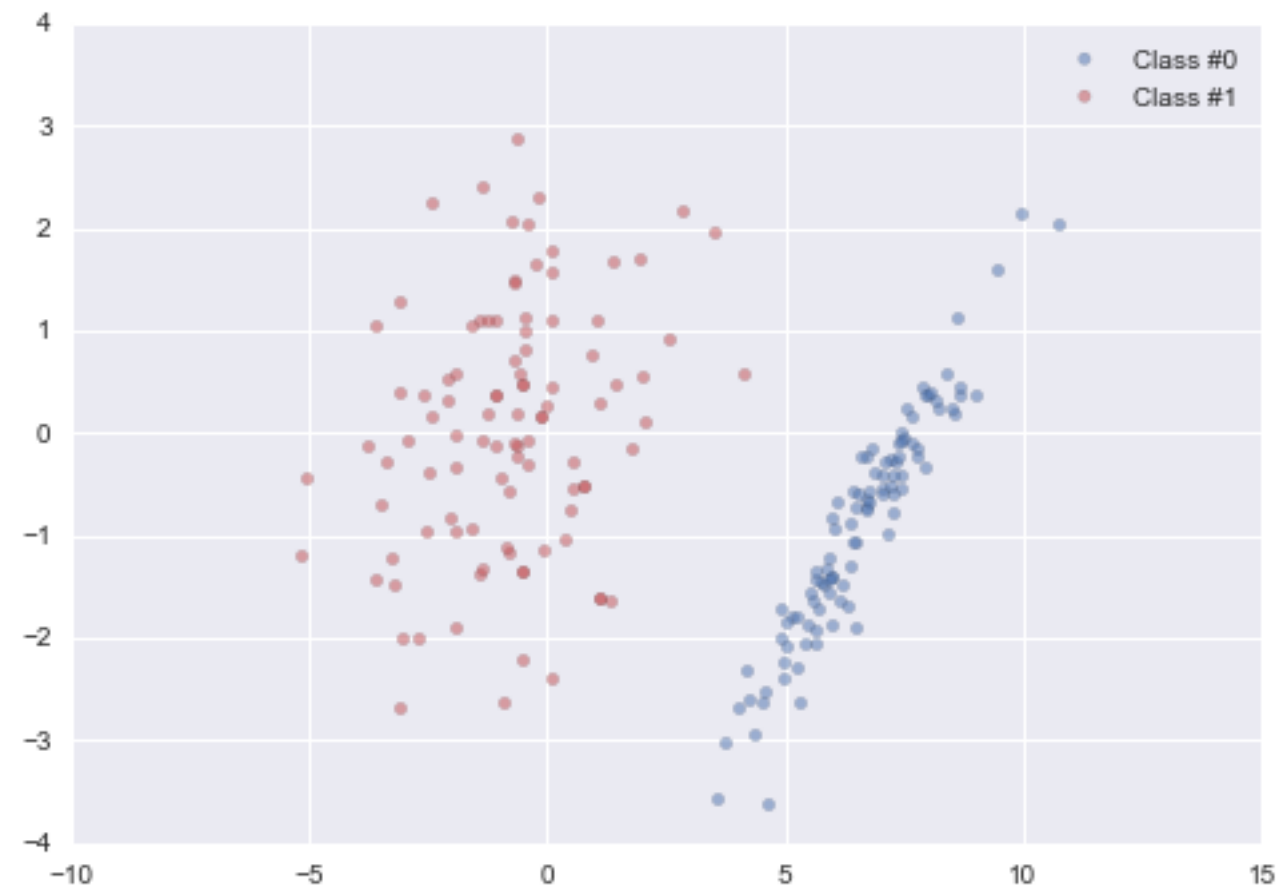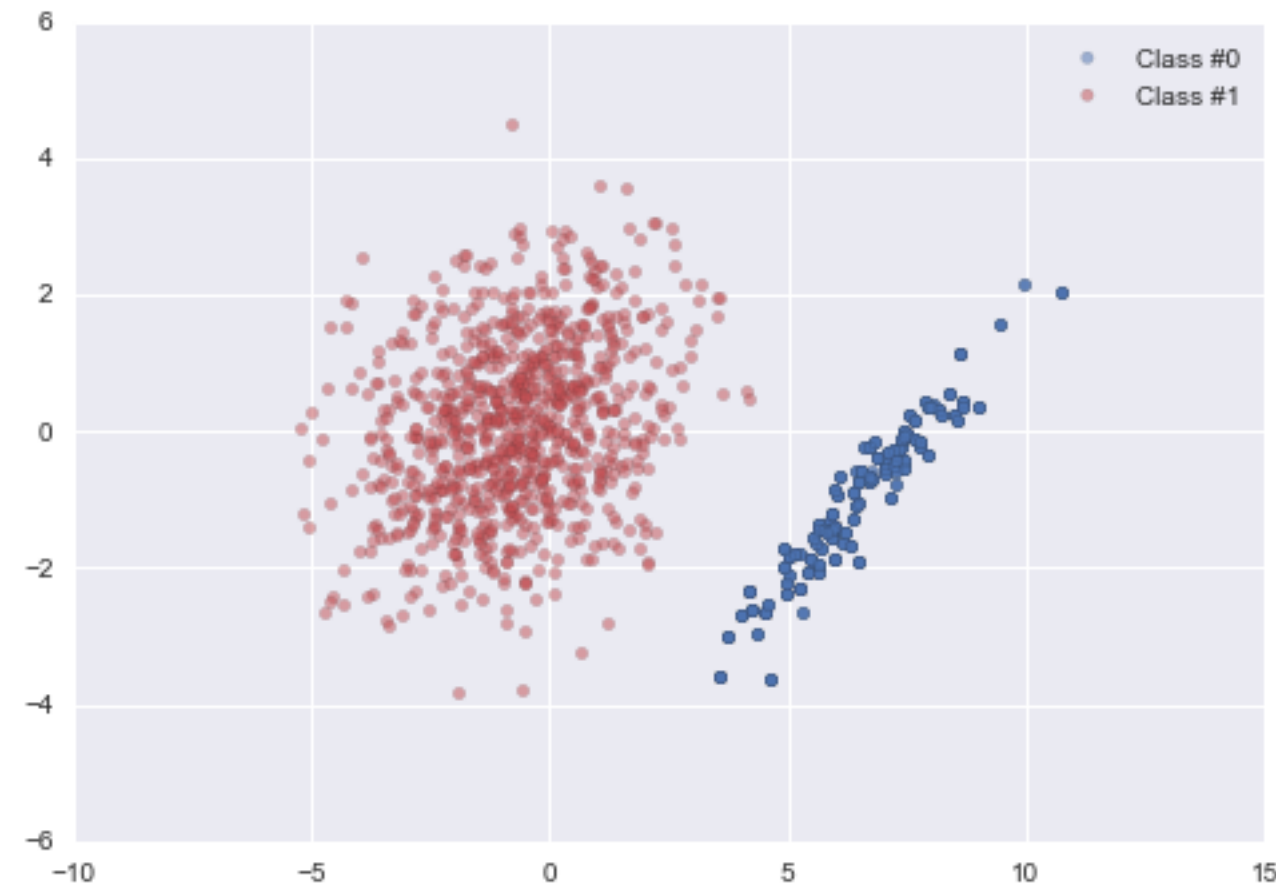# Undersampling

- Randomly discards majority class observations

- Pros: Makes calculations way faster

- Cons: Throwing out data :(

# Oversampling

- Replicates minority class observations

- Pros: Doesn't discard info
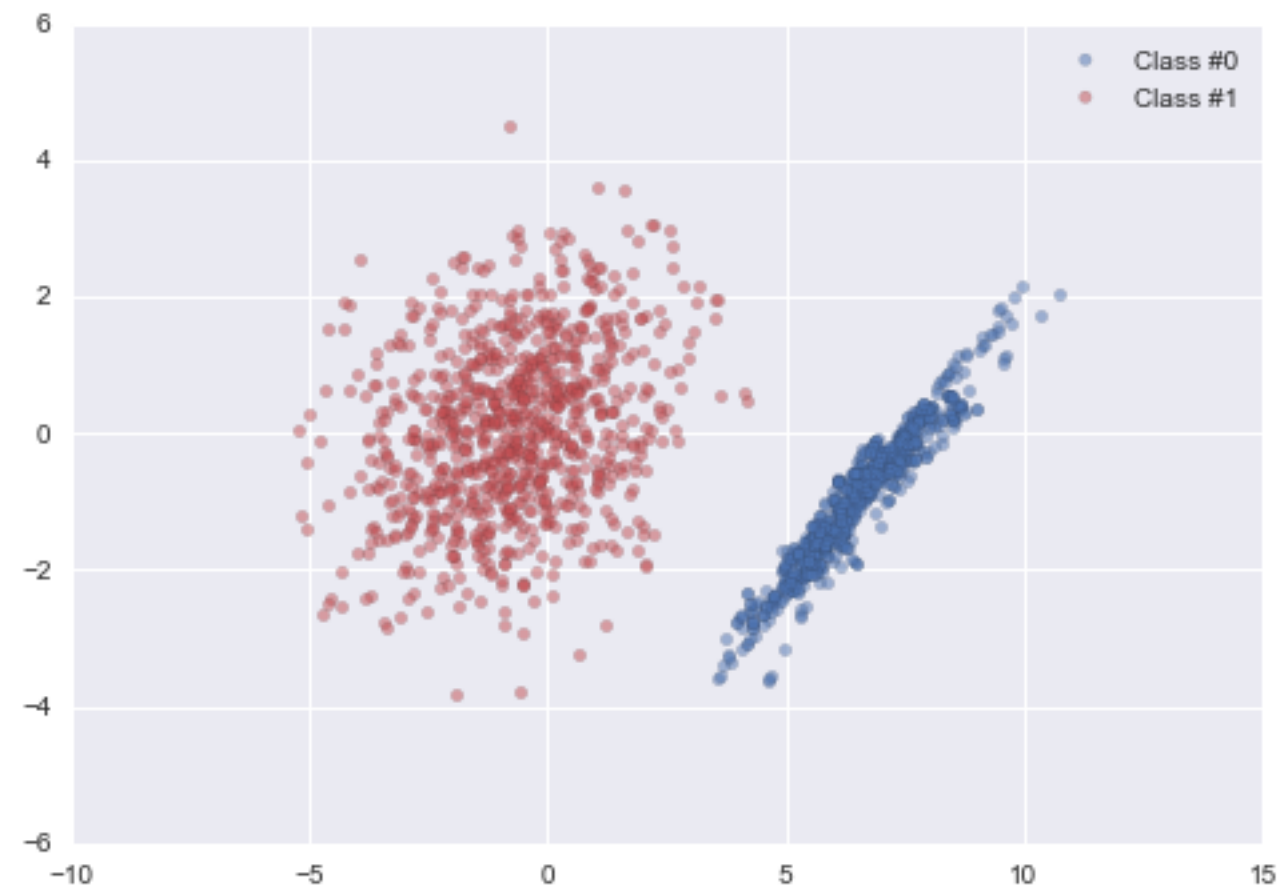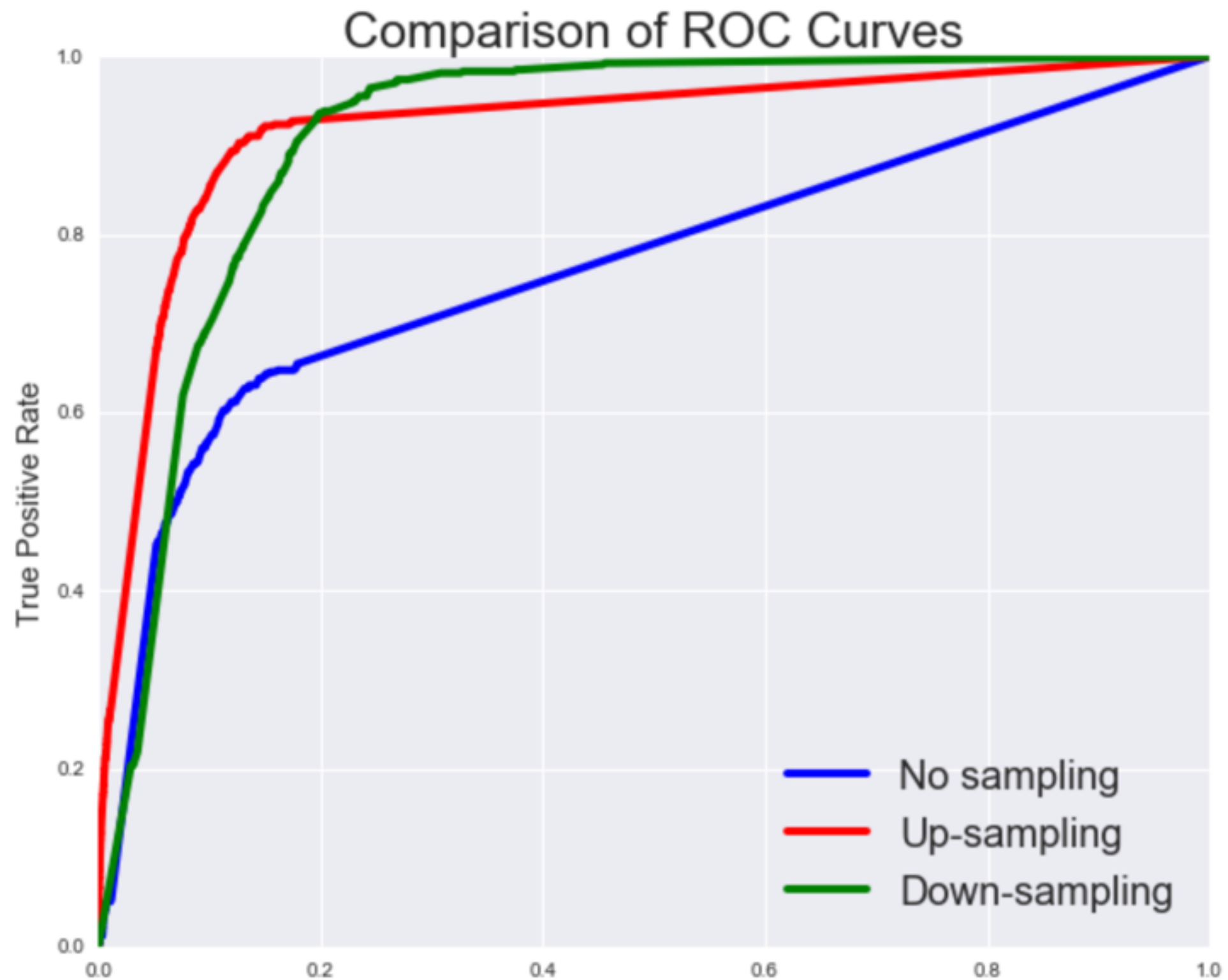
- Cons: Overfitting likely :(

# SMOTE

- Generates synthetic minority class observations Let's look at the original paper's pseudocode:

  https://www.jair.org/media/953/live-953-2037-jair.pdf

# How do we pick?



Comparison of ROC Curves

True Positive Rate

— No sampling
— Up-sampling
— Down-sampling

# Changing Cost Function

- Models with explicit cost function can be modified to incorporate classification cost

  - e.g. SVM, logistic $\qquad \dfrac{\|w\|^2}{2} + C\sum_{i=1}^{n}\xi_i \longrightarrow \dfrac{\|w\|^2}{2} + C^+ \sum_{\{i|yi=+1\}}^{n_+}\xi_i + C^- \sum_{\{j|yj=-1\}}^{n_-}\xi_j$

- Can affect the optimization

  - ex. cost sensitive logistic regression no longer convex

- Not possible for all models