

Carrying Out Error Analysis

< Look at dev examples to evaluate ideas >

Should you try to make your cat classifier do better on dogs? focus on the dog problem?

90% accuracy
10% error

Error analysis: 5-10 min

- Get ~100 mislabeled dev set examples. $\downarrow 5\%$ 5/100
- Count up how many are dogs. $\downarrow 10\%$ error decrease. 9.5% error

$\downarrow 70\%$ 70/100 $\downarrow 10\%$ effective $\downarrow 5\%$

"ceiling"
what's in the best case how well could working on the dog problem help you

"error analysis": effective to decide what is the most important thing / what to focus.

< Evaluate multiple ideas in parallel >

Ideas for cat detection:

- Fix picture or dogs being recognized as cats
- Fix great cats (lions, panthers, etc...) being misrecognized
- Improve performance on blurry images.

draw the table

estimate of how worthwhile it might be to work on each of these different categories of errors.

Image	Dog	Great Cats	Blurry	Insta gram	Comments
1	✓			✓	Pitbull
2					
3				✓	
...					
% of total	8%	43%	61%	12%	Rainy day at zoo

ceiling of data is much higher.

Cleaning up Incorrectly labeled data

Incorrectly labeled data \rightarrow fix it!

< Incorrectly labeled examples >

x	cat	dog	cat	cat	dog	dog	cat
y	1	0	1	1	0	1	1

incorrect label.
 \downarrow Should be 0 but get wrong.

Training set.

DL Algorithms are quite robust to random errors in the training set.

Systematic errors: less robust.

< Error Analysis >

Image	Dog	Great Cat	Blurry	Incorrectly labeled	Comments
...					
q8				✓	Labeler missed cat in background
q9		✓			
100				✓	Drawing of a cat; Not a real cat.
% of total	8%	43%	61%	6%	
Overall dev set error		10%			
Errors due incorrect labels		0.6%			2% higher fraction
Errors due to other causes		9.4%			1.4% fix ok

fixing it is not the most incorrect right now (0.6% - very small)

focus on fixing!

fixing it is not the most incorrect right now (0.6% - very small)


Goal of dev set is to help you select between two classifiers A & B.

if it makes a significant difference to your ability to evaluate algorithms on your dev set, then go ahead and spend the time to fix incorrect labels. doesn't make a significant difference \rightarrow not be the best use of your time.

<correcting incorrect dev/test set examples>

- Apply same process to your dev and test sets to make sure they continue to come from the same distribution.
- Consider examining examples your algorithm got right as well }
as ones it got wrong. 7% 98%
↳ doesn't always done. Just something to consider
- Train and dev/test data may now come from slightly different distributions.

Build First System Quickly, Then Iterate

- Set up dev / test set and metric  set up a target
- build initial system quickly
- Use bias / variance analysis & Error analysis to prioritize next steps.