



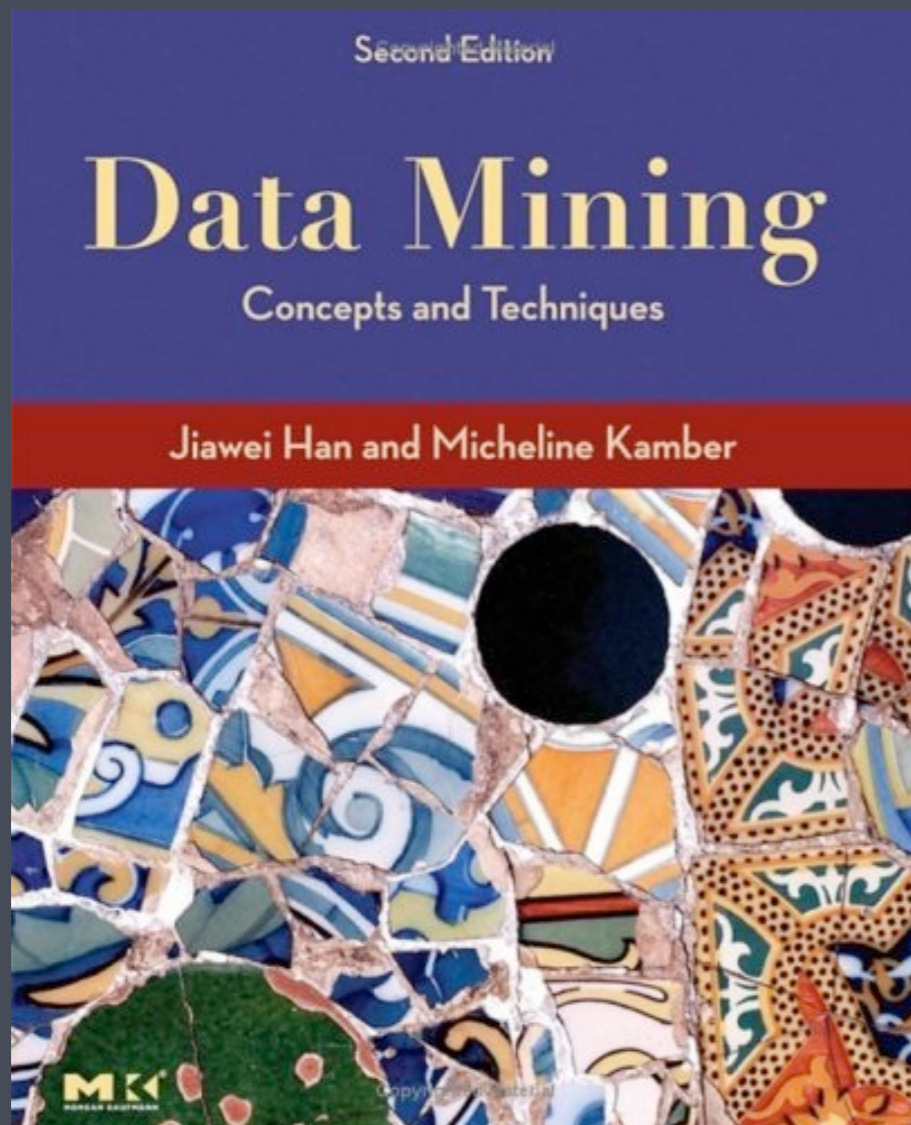
MINING SOFTWARE REPOSITORIES

Software Engineering Course 2009

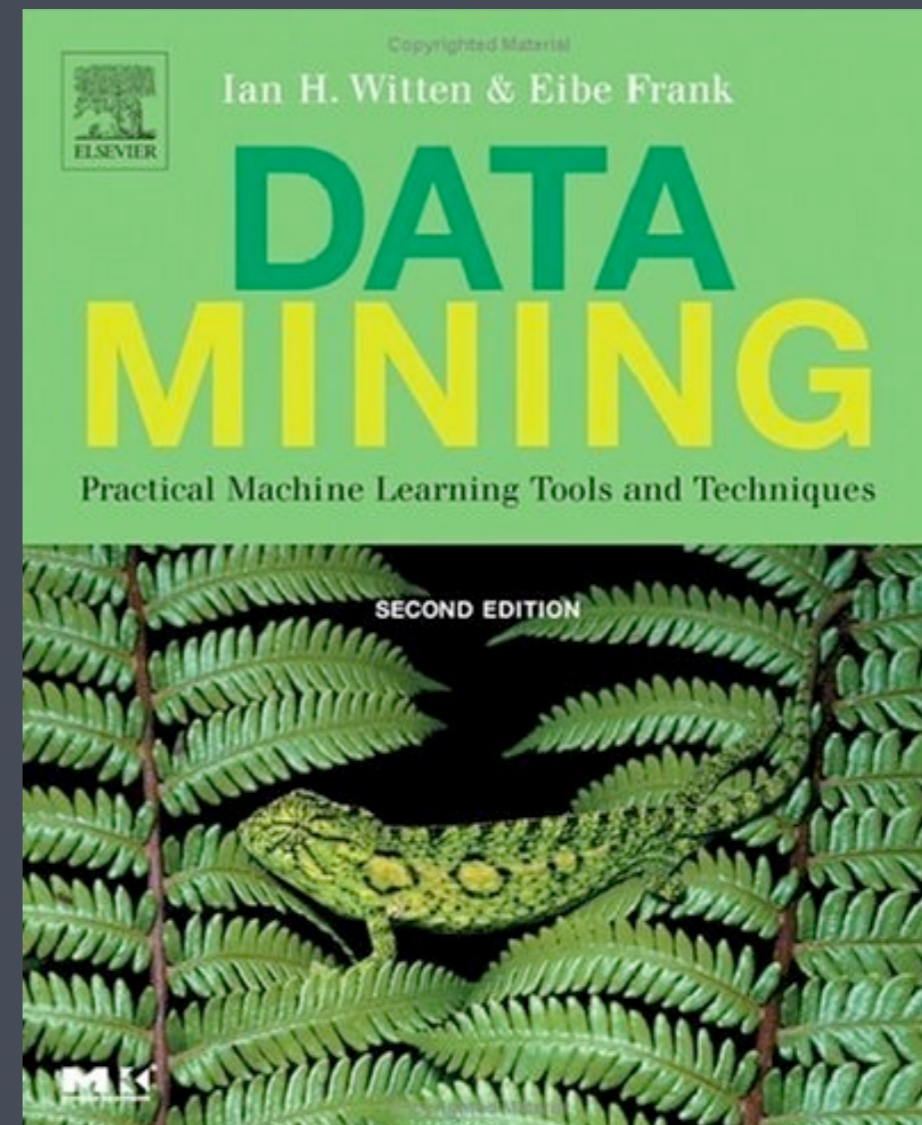
Kim Herzig - Saarland University



Books



Data Mining: Concepts and Techniques
by **Jiawei Han & Micheline Kamber**



Data Mining: Practical Machine Learning Tools and Techniques
by **Ian H. Witten & Eibe Frank**



Imagine



You as Quality Manager



Imagine



- ▶ 30,000 classes
- ▶ ~ 5.5 million lines of code
- ▶ ~3000 defect per release
- ▶ 700 developers

You as Quality Manager

Your product



Your Boss

Test the system!
You have 6 months, \$500,000.
And don't miss any bug!





The Problem

- ▶ Not enough **time** to test everything
 - ▶ What to test? What to test first?
- ▶ Not enough **money** to pay enough testers
 - ▶ To which extend?

Central question:

Where are the most defect prone entities in my system?

Your Testers

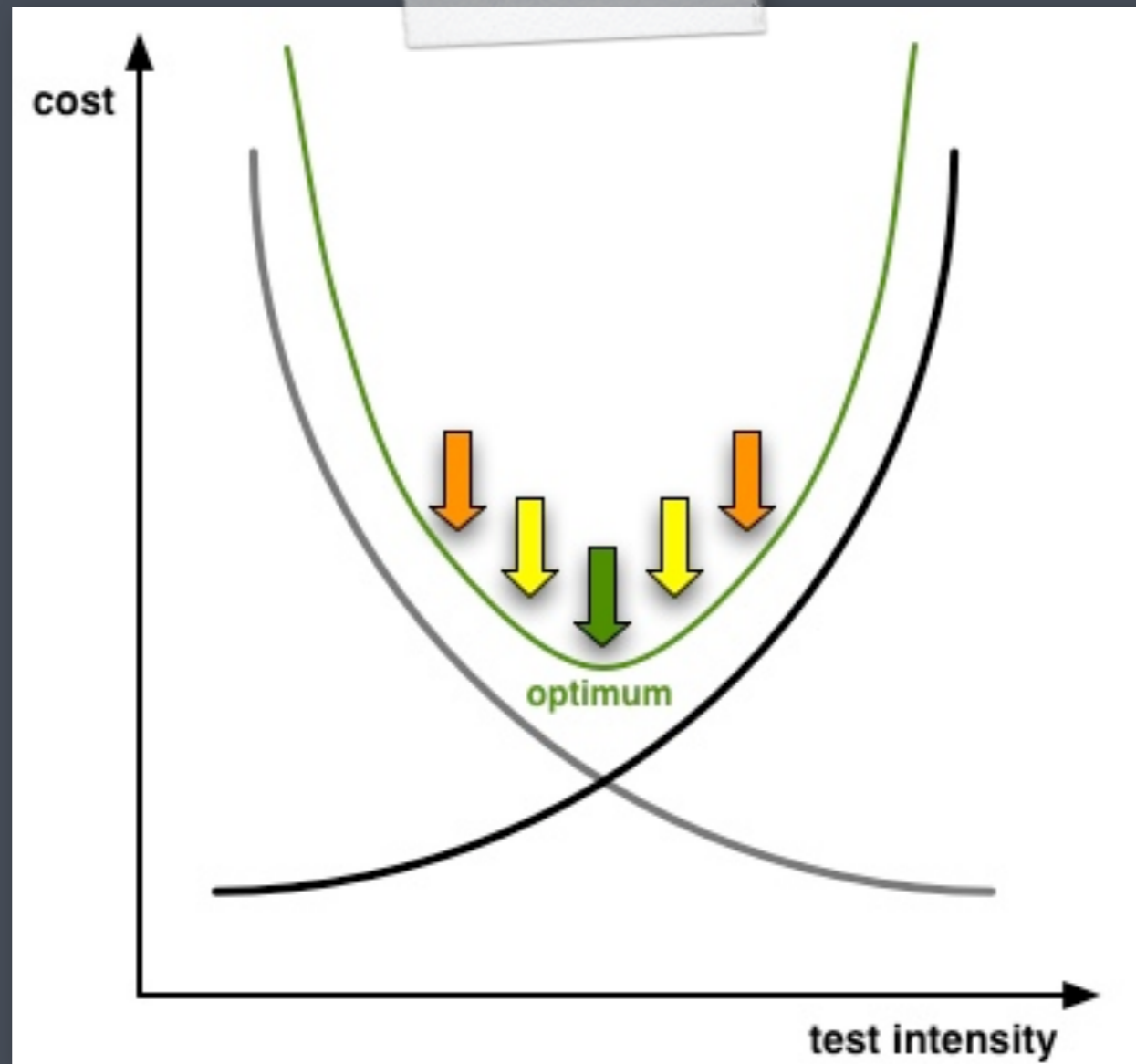
Your Testers



We need efficiency!

We need efficiency!

We need efficiency!



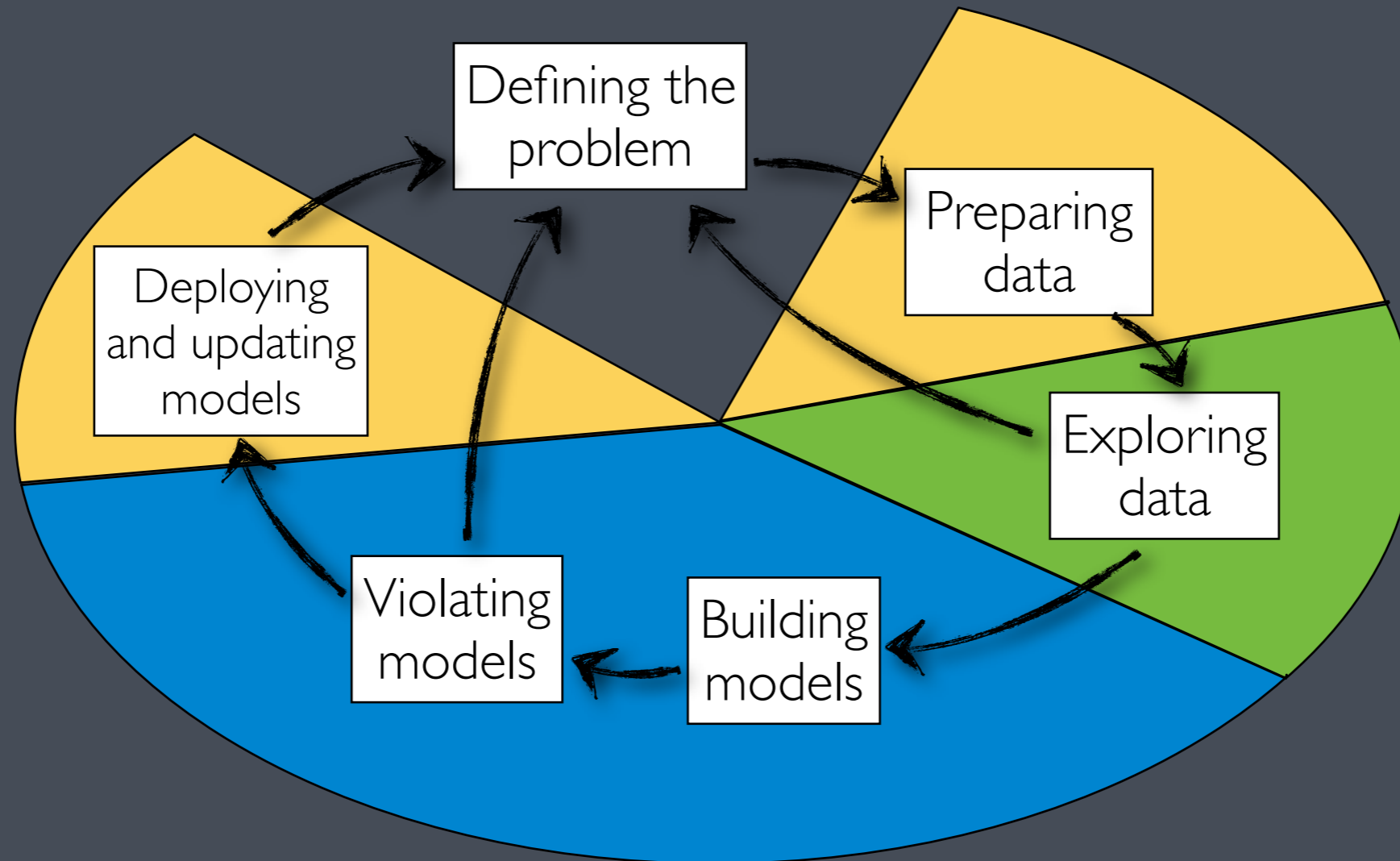
Can we learn from history?
... to predict or estimate the future?

data mining

What is data mining?

Data mining is the process of discovering actionable information from large sets of data.

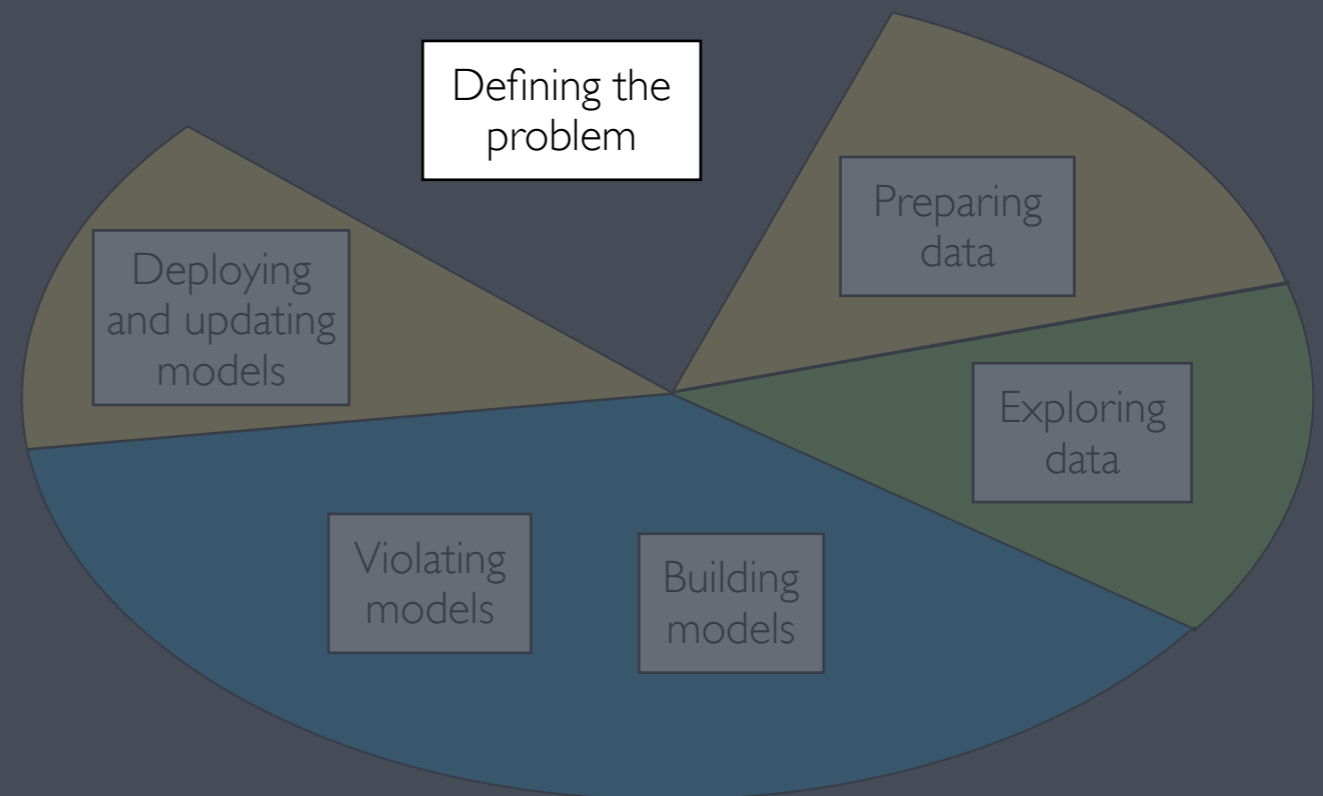
The Mining Model



<http://technet.microsoft.com/en-us/library/ms174949.aspx>

Step 1: Defining Problem

- ▶ Clearly define the problem
 - ▶ What are you looking for?
 - ▶ Scope of problem
 - ▶ Types of relationships
- ▶ Define how to evaluate
 - ▶ Prediction, recommendation or just patterns



Defect Prediction Problem

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model



Which source code entities
should we test most?

Defect Prediction Problem

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model

Which source code entities should we test most?

Which are the most defect prone entities in my system?



Defect Prediction Problem

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model

Which source code entities should we test most?

Which are the most defect prone entities in my system?

In the past, which entities had the most defects?



Defect Prediction Problem

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model

Which properties of source code entities correlate with defects?

Which source code entities should we test most?

Which are the most defect prone entities in my system?

In the past, which entities had the most defects?

Data Sources

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model



Bug Database



Version Archive



Source Code



Data Sources

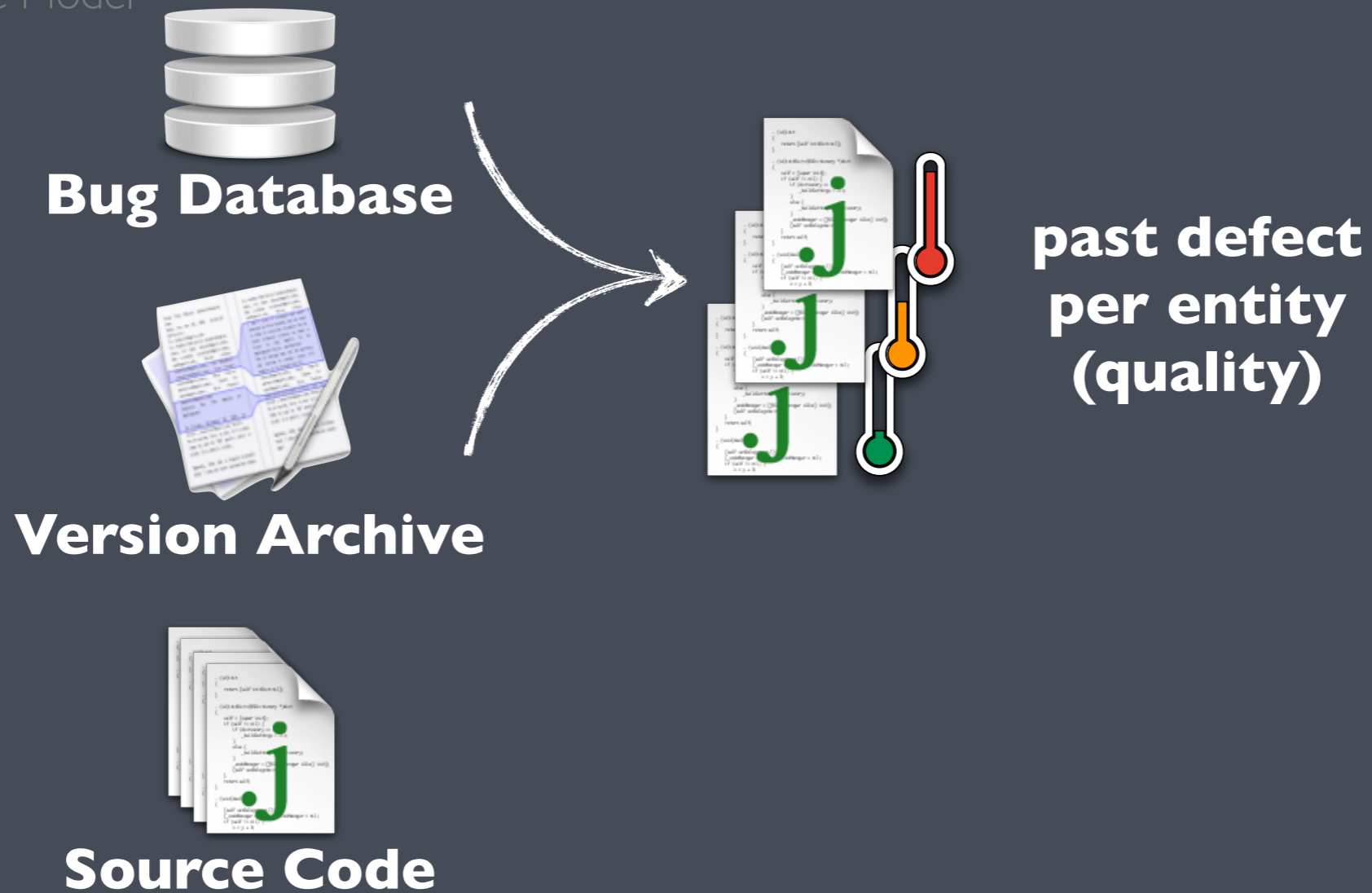
Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model



Data Sources

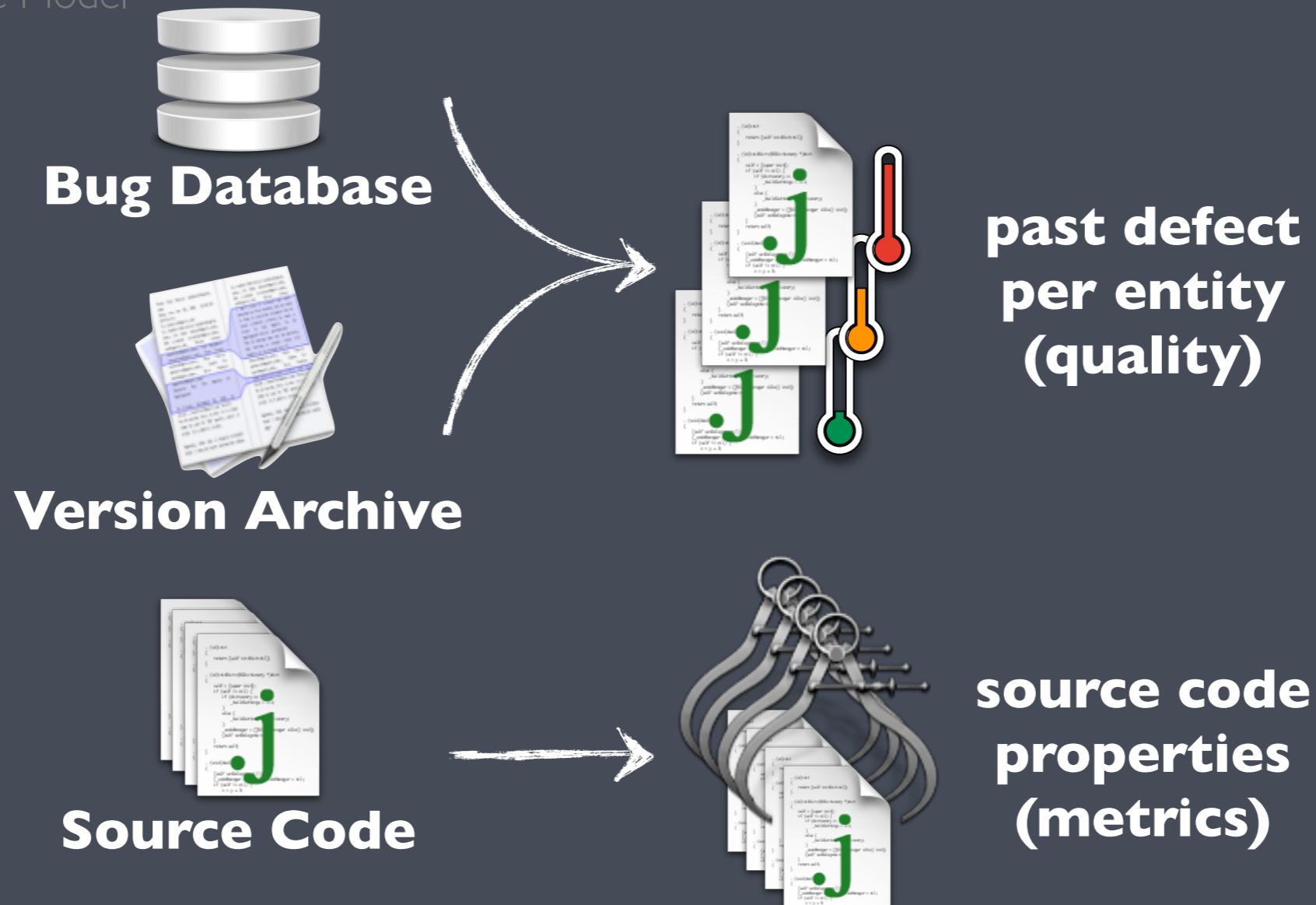
Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model



Data Sources: Heuristics

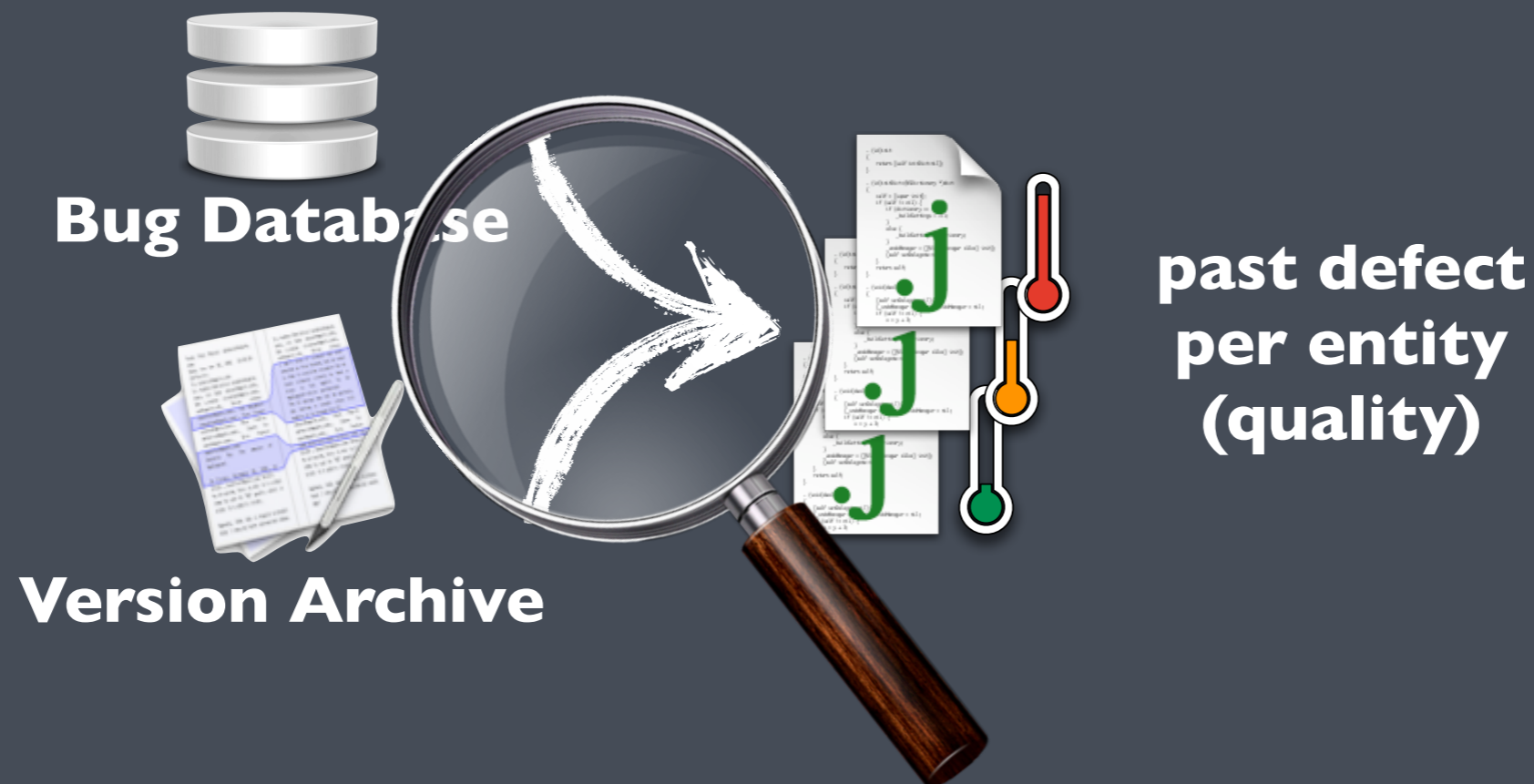
Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model



"... commit messages that contain fix and bug id ..."

Data Sources: Metrics

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

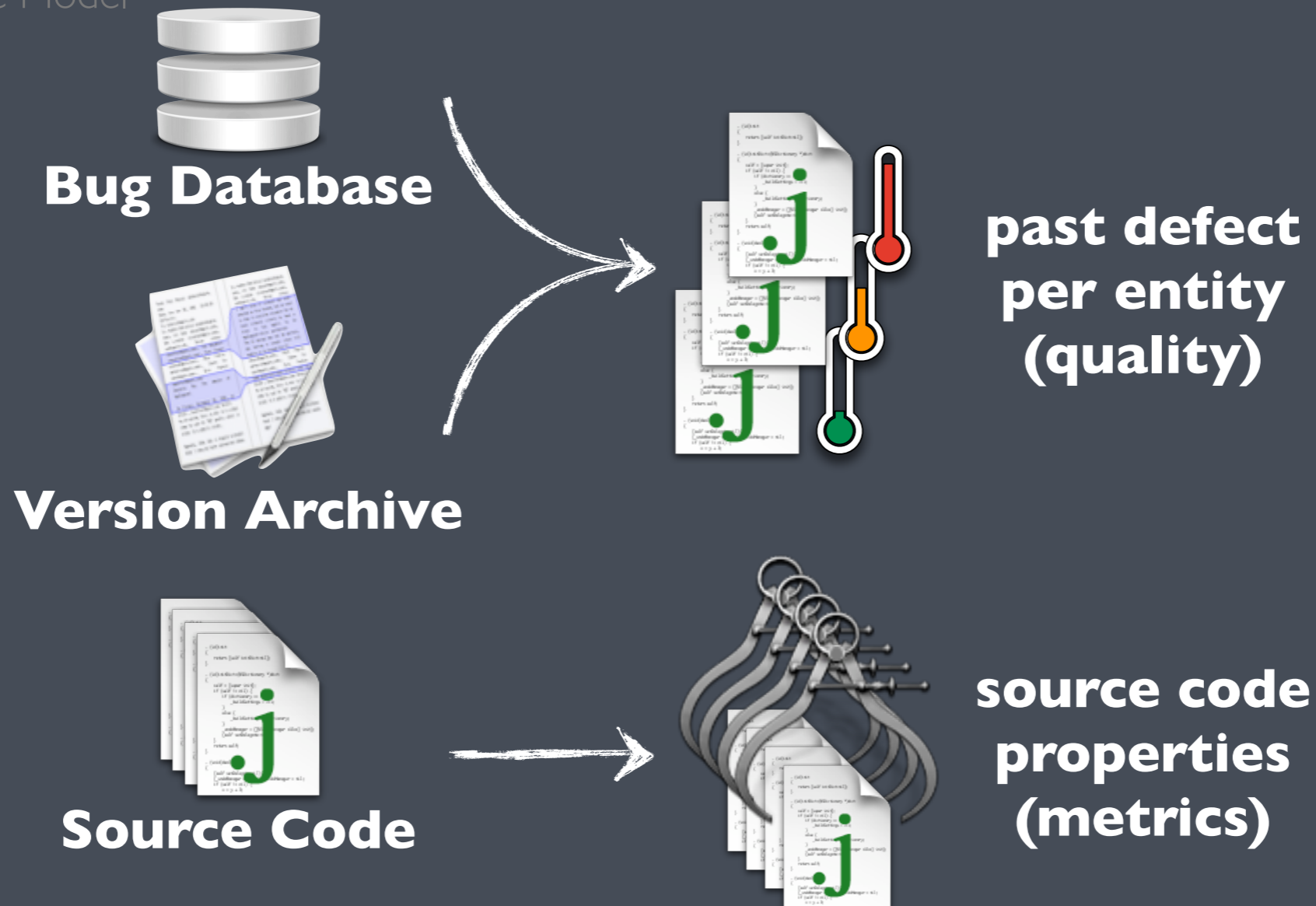
Step 5: Validating the Model

- ▶ Complicity metrics
 - ▶ McCabe, FanIn, FanOut, Couplings
 - ▶ (see Lecture “Metrics and Estimation”)
- ▶ Time metrics
 - ▶ How many changes
 - ▶ How many different authors
 - ▶ Age of code



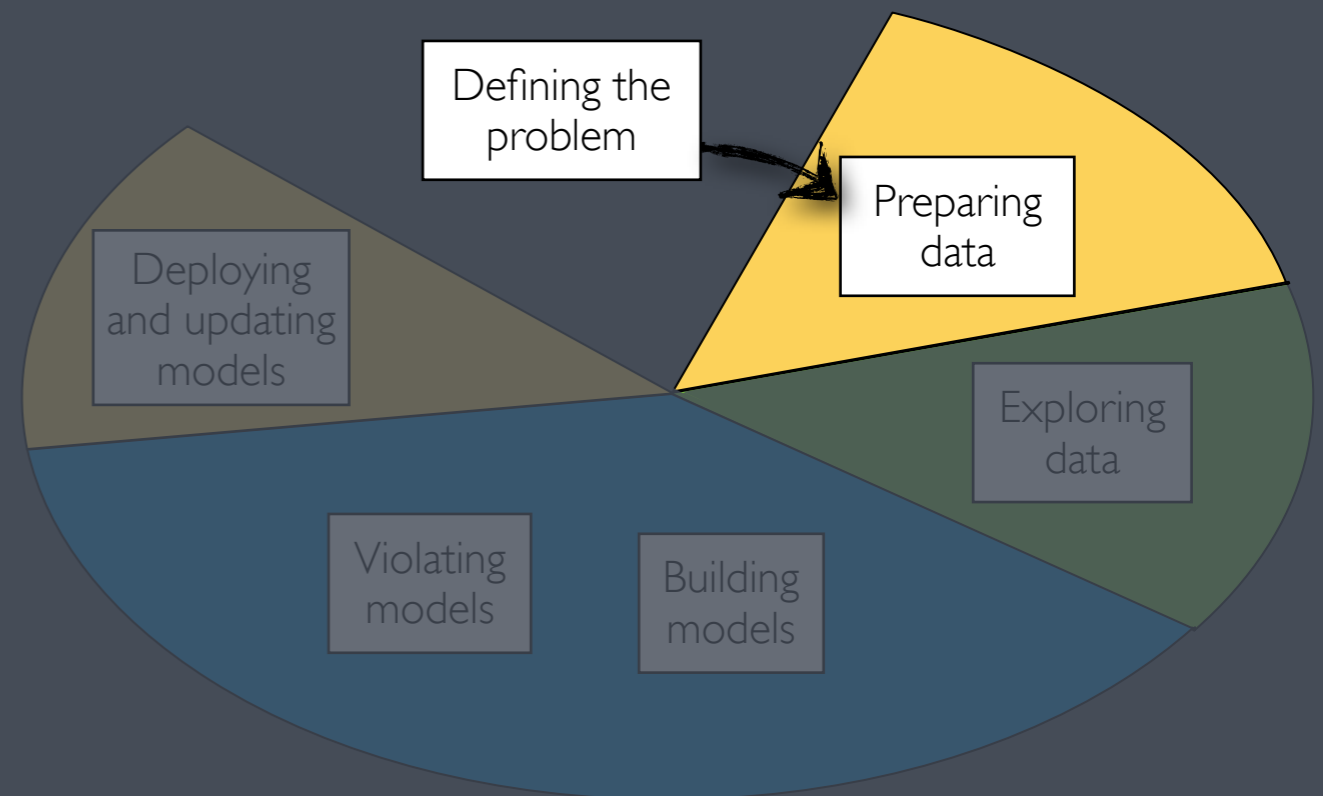
Data Sources

- Step 1: Define the problem
- Step 2: Prepare Data**
- Step 3: Explore Data
- Step 4: Building the Model
- Step 5: Validating the Model



Step 2: Prepare Data

- ▶ Highly distributed data:
 - ▶ Version repository, bug data base, time trackers, ...
- ▶ Data integration
 - ▶ Excel, CSV, SQL, ARFF, ...
- ▶ Data cleaning
 - ▶ missing values, noise, inter-correlations



Example Mining File

Package	Name	bugs	No. Methods	CCV	LCOM	TCC	MAXCC	AVCC	NOS	HEFF	UWCS
org.jruby.compiler	ASTCompiler	0	278	605	1	605	92	2.18	1574	424794.18	279
org.jruby.compiler	ASTCompiler\$OpElementAsgnArgumentsCallback	0	3	6	0	6	3	2	21	3592.03	4
org.jruby.compiler	ASTCompiler\$SpecificArityArguments	6	3	7	0.5	7	4	2.33	22	2660.26	5
org.jruby.compiler	ASTCompiler\$VariableArityArguments	0	3	3	0.5	3	1	1	7	31.61	4
org.jruby.compiler	ASTCompiler19	1	20	40	0	40	10	2	117	26264.39	20
org.jruby.compiler	ASTInspector	0	18	51	1.05	51	25	2.83	433	562739.87	45
org.jruby.evaluator	ASTInterpreter	4	22	55	0	55	6	2.5	169	38462.7	22
org.jruby.runtime	AbstractCompiledBlockCallback	0	0	0	0	0	0	0	1	0	0
org.jruby.ext.ffi	AbstractInvoker	5	7	11	0.92	11	4	1.57	42	6346.91	9
org.jruby.ext.ffi	AbstractMemory	0	68	105	0.97	105	5	1.54	310	87618.03	71
org.jruby.ast.executable	AbstractScript	1	166	187	24	187	3	1.13	413	28142.16	190
org.jruby.compiler.impl	AbstractVariableCompiler	8	20	61	0.71	61	12	3.05	279	129005.27	27
org.jruby.util	Adler32Ext	0	9	10	0.67	10	2	1.11	27	484	12
org.jruby.internal.runtime.methods	AliasMethod	2	40	45	0.96	45	6	1.12	102	1383.47	53
org.jruby.ast	AliasNode	2	7	7	0.67	7	1	1	24	90.68	9
org.jruby.ext.ffi	AllocatedDirectMemoryIO	11	2	2	0	2	1	1	3	0	2
org.jruby.ext.ffi.jffi	AllocatedNativeMemoryIO	3	5	8	0.75	8	2	1.6	20	1138.3	8
org.jruby.ext.ffi.jna	AllocatedNativeMemoryIO	0	4	6	0.67	6	2	1.5	19	884.91	5
org.jruby.ast	AndNode	2	7	8	0.5	8	2	1.14	27	871.77	9
org.jruby.anno	AnnotationBinder	6	3	3	1	3	1	1	27	32.73	5
org.jruby.anno	AnnotationBinder\$AnnotationBindingProcessor	2	5	14	0.88	14	6	2.8	41	2180.79	7
org.jruby.anno	AnnotationBinder\$AnnotationBindingProcessor\$RubyClassVisitor	0	8	66	0.86	66	22	8.25	246	230539.24	10
org.jruby.ast	ArgAuxillaryNode	8	6	6	0.8	6	1	1	17	85.47	8
org.jruby.ast	ArgsCatNode	2	7	7	0.5	7	1	1	26	806.1	9
org.jruby.ast	ArgsNoArgNode	10	2	2	0	2	1	1	8	0	2
org.jruby.ast	ArgsNode	4	34	64	0.9	64	7	1.88	146	33903.67	48
org.jruby.ast	ArgsPreOneArgNode	11	6	6	0	6	1	1	17	115.79	6
org.jruby.ast	ArgsPreTwoArgNode	2	7	8	0	8	2	1.14	20	231.46	7
org.jruby.ast	ArgsPushNode	2	7	7	0.67	7	1	1	26	642.14	9
org.jruby.ast.util	ArgsUtil	0	5	13	0	13	3	2.6	35	4082.78	5
org.jruby.ast	ArgumentNode	1	6	6	0.6	6	1	1	17	53.48	7
org.jruby.compiler	ArgumentsCallback	1	1	1	0	1	1	1	2	0	1
org.jruby.runtime	Arity	12	26	55	1.01	55	7	2.12	162	34258.83	37

Example Mining File

Package	Name	bugs	No. Methods	CCV	LCOM	TCC	MAXCC	AVCC	NOS	HEFF	UWCS
org.jruby.compiler	ASTCompiler	0	278	605	1	605	92	2.18	1574	424794.18	279
org.jruby.compiler	ASTCompiler\$OpElementAsgnArgumentsCallback	0	3	6	0	6	3	2	21	3592.03	4
org.jruby.compiler	ASTCompiler\$SpecificArityArguments	6	3	7	0.5	7	4	2.33	22	2660.26	5
org.jruby.compiler	ASTCompiler\$VariableArityArguments	0	3	3	0.5	3	1	1	7	31.61	4
org.jruby.compiler	ASTCompiler19	1	20	40	0	40	10	2	117	26264.39	20
org.jruby.compiler	ASTInspector	0	18	51	1.05	51	25	2.83	433	562739.87	45
org.jruby.evaluator	ASTInterpreter	4	22	55	0	55	6	2.5	169	38462.7	22
org.jruby.runtime	AbstractCompiledBlockCallback	0	0	0	0	0	0	0	1	0	0
org.jruby.ext.ffi	AbstractInvoker	5	7	11	0.92	11	4	1.57	42	6346.91	9
org.jruby.ext.ffi	AbstractMemory	0	68	105	0.97	105	5	1.54	310	87618.03	71
org.jruby.ast.executable	AbstractScript	1	166	187	24	187	3	1.13	413	28142.16	190
org.jruby.compiler.impl	AbstractVariableCompiler	8	20	61	0.71	61	12	3.05	279	129005.27	27
org.jruby.util	Adler32Ext	0	9	10	0.67	10	2	1.11	27	484	12
org.jruby.internal.runtime.methods	AliasMethod	2	40	45	0.96	45	6	1.12	102	1383.47	53
org.jruby.ast	AliasNode	2	7	7	0.67	7	1	1	24	90.68	9
org.jruby.ext.ffi	AllocatedDirectMemoryIO	11	2	2	0	2	1	1	3	0	2
org.jruby.ext.ffi.jffi	AllocatedNativeMemoryIO	3	5	8	0.75	8	2	1.6	20	1138.3	8
org.jruby.ext.ffi.jna	AllocatedNativeMemoryIO	0	4	6	0.67	6	2	1.5	19	884.91	5
org.jruby.ast	AndNode	2	7	8	0.5	8	2	1.14	27	871.77	9
org.jruby.anno	AnnotationBinder	6	3	3	1	3	1	1	27	32.73	5
org.jruby.anno	AnnotationBinder\$AnnotationBindingProcessor	2	5	14	0.88	14	6	2.8	41	2180.79	7
org.jruby.anno	AnnotationBinder\$AnnotationBindingProcessor\$RubyClassVisitor	0	8	66	0.86	66	22	8.25	246	230539.24	10
org.jruby.ast	ArgAuxillaryNode	8	6	6	0.8	6	1	1	17	85.47	8
org.jruby.ast	ArgsCatNode	2	7	7	0.5	7	1	1	26	806.1	9
org.jruby.ast	ArgsNoArgNode	10	2	2	0	2	1	1	8	0	2
org.jruby.ast	ArgsNode	4	34	64	0.9	64	7	1.88	146	33903.67	48
org.jruby.ast	ArgsPreOneArgNode	11	6	6	0	6	1	1	17	115.79	6
org.jruby.ast	ArgsPreTwoArgNode	2	7	8	0	8	2	1.14	20	231.46	7
org.jruby.ast	ArgsPushNode	2	7	7	0.67	7	1	1	26	642.14	9
org.jruby.ast.util	ArgsUtil	0	5	13	0	13	3	2.6	35	4082.78	5
org.jruby.ast	ArgumentNode	1	6	6	0.6	6	1	1	17	53.48	7
org.jruby.compiler	ArgumentsCallback	1	1	1	0	1	1	1	2	0	1
org.jruby.runtime	Arity	12	26	55	1.01	55	7	2.12	162	34258.83	37

entities

Example Mining File

Package	Name	bugs	No. Methods	CCV	LCOM	TCC	MAXCC	AVCC	NOS	HEFF	UWCS
org.jruby.compiler	ASTCompiler	0	278	605	1	605	92	2.18	1574	424794.18	279
org.jruby.compiler	ASTCompiler\$OpElementAsgnArgumentsCallback	0	3	6	0	6	3	2	21	3592.03	4
org.jruby.compiler	ASTCompiler\$SpecificArityArguments	6	3	7	0.5	7	4	2.33	22	2660.26	5
org.jruby.compiler	ASTCompiler\$VariableArityArguments	0	3	3	0.5	3	1	1	7	31.61	4
org.jruby.compiler	ASTCompiler19	1	20	40	0	40	10	2	117	26264.39	20
org.jruby.compiler	ASTInspector	0	18	51	1.05	51	25	2.83	433	562739.87	45
org.jruby.evaluator	ASTInterpreter	4	22	55	0	55	6	2.5	169	38462.7	22
org.jruby.runtime	AbstractCompiledBlockCallback	0	0	0	0	0	0	0	1	0	0
org.jruby.ext.ffi	AbstractInvoker	5	7	11	0.92	11	4	1.57	42	6346.91	9
org.jruby.ext.ffi	AbstractMemory	0	68	105	0.97	105	5	1.54	310	87618.03	71
org.jruby.ast.executable	AbstractScript	1	166	187	24	187	3	1.13	413	28142.16	190
org.jruby.compiler.impl	AbstractVariableCompiler	8	20	61	0.71	61	12	3.05	279	129005.27	27
org.jruby.util	Adler32Ext	0	9	10	0.67	10	2	1.11	27	484	12
org.jruby.internal.runtime.methods	AliasMethod	2	40	45	0.96	45	6	1.12	102	1383.47	53
org.jruby.ast	AliasNode	2	7	7	0.67	7	1	1	24	90.68	9
org.jruby.ext.ffi	AllocatedDirectMemoryIO	11	2	2	0	2	1	1	3	0	2
org.jruby.ext.ffi.jffi	AllocatedNativeMemoryIO	3	5	8	0.75	8	2	1.6	20	1138.3	8
org.jruby.ext.ffi.jna	AllocatedNativeMemoryIO	0	4	6	0.67	6	2	1.5	19	884.91	5
org.jruby.ast	AndNode	2	7	8	0.5	8	2	1.14	27	871.77	9
org.jruby.anno	AnnotationBinder	6	3	3	1	3	1	1	27	32.73	5
org.jruby.anno	AnnotationBinder\$AnnotationBindingProcessor	2	5	14	0.88	14	6	2.8	41	2180.79	7
org.jruby.anno	AnnotationBinder\$AnnotationBindingProcessor\$RubyClassVisitor	0	8	66	0.86	66	22	8.25	246	230539.24	10
org.jruby.ast	ArgAuxillaryNode	8	6	6	0.8	6	1	1	17	85.47	8
org.jruby.ast	ArgsCatNode	2	7	7	0.5	7	1	1	26	806.1	9
org.jruby.ast	ArgsNoArgNode	10	2	2	0	2	1	1	8	0	2
org.jruby.ast	ArgsNode	4	34	64	0.9	64	7	1.88	146	33903.67	48
org.jruby.ast	ArgsPreOneArgNode	11	6	6	0	6	1	1	17	115.79	6
org.jruby.ast	ArgsPreTwoArgNode	2	7	8	0	8	2	1.14	20	231.46	7
org.jruby.ast	ArgsPushNode	2	7	7	0.67	7	1	1	26	642.14	9
org.jruby.ast.util	ArgsUtil	0	5	13	0	13	3	2.6	35	4082.78	5
org.jruby.ast	ArgumentNode	1	6	6	0.6	6	1	1	17	53.48	7
org.jruby.compiler	ArgumentsCallback	1	1	1	0	1	1	1	2	0	1
org.jruby.runtime	Arity	12	26	55	1.01	55	7	2.12	162	34258.83	37

entities

data points



Example Mining File

Package	Name	bugs	No. Methods	CCV	LCOM	TCC	MAXCC	AVCC	NOS	HEFF	UWCS
org.jruby.compiler	ASTCompiler	0	278	605	1	605	92	2.18	1574	424794.18	279
org.jruby.compiler	ASTCompiler\$OpElementAsgnArgumentsCallback	0	3	6	0	6	3	2	21	3592.03	4
org.jruby.compiler	ASTCompiler\$SpecificArityArguments	6	3	7	0.5	7	4	2.33	22	2660.26	5
org.jruby.compiler	ASTCompiler\$VariableArityArguments	0	3	3	0.5	3	1	1	7	31.61	4
org.jruby.compiler	ASTCompiler19	1	20	40	0	40	10	2	117	26264.39	20
org.jruby.compiler	ASTInspector	0	18	51	1.05	51	25	2.83	433	562739.87	45
org.jruby.evaluator	ASTInterpreter	4	22	55	0	55	6	2.5	169	38462.7	22
org.jruby.runtime	AbstractCompiledBlockCallback	0	0	0	0	0	0	0	1	0	0
org.jruby.ext.ffi	AbstractInvoker	5	7	11	0.92	11	4	1.57	42	6346.91	9
org.jruby.ext.ffi	AbstractMemory	0	68	105	0.97	105	5	1.54	310	87618.03	71
org.jruby.ast.executable	AbstractScript	1	166	187	24	187	3	1.13	413	28142.16	190
org.jruby.compiler.impl	AbstractVariableCompiler	8	20	61	0.71	61	12	3.05	279	129005.27	27
org.jruby.util	Adler32Ext	0	9	10	0.67	10	2	1.11	27	484	12
org.jruby.internal.runtime.methods	AliasMethod	2	40	45	0.96	45	6	1.12	102	1383.47	53
org.jruby.ast	AliasNode	2	7	7	0.67	7	1	1	24	90.68	9
org.jruby.ext.ffi	AllocatedDirectMemoryIO	11	2	2	0	2	1	1	3	0	2
org.jruby.ext.ffi.jffi	AllocatedNativeMemoryIO	3	5	8	0.75	8	2	1.6	20	1138.3	8
org.jruby.ext.ffi.jna	AllocatedNativeMemoryIO	0	4	6	0.67	6	2	1.5	19	884.91	5
org.jruby.ast	AndNode	2	7	8	0.5	8	2	1.14	27	871.77	9
org.jruby.anno	AnnotationBinder	6	3	3	1	3	1	1	27	32.73	5
org.jruby.anno	AnnotationBinder\$AnnotationBindingProcessor	2	5	14	0.88	14	6	2.8	41	2180.79	7
org.jruby.anno	AnnotationBinder\$AnnotationBindingProcessor\$RubyClassVisitor	0	8	66	0.86	66	22	8.25	246	230539.24	10
org.jruby.ast	ArgAuxillaryNode	8	6	6	0.8	6	1	1	17	85.47	8
org.jruby.ast	ArgsCatNode	2	7	7	0.5	7	1	1	26	806.1	9
org.jruby.ast	ArgsNoArgNode	10	2	2	0	2	1	1	8	0	2
org.jruby.ast	ArgsNode	4	34	64	0.9	64	7	1.88	146	33903.67	48
org.jruby.ast	ArgsPreOneArgNode	11	6	6	0	6	1	1	17	115.79	6
org.jruby.ast	ArgsPreTwoArgNode	2	7	8	0	8	2	1.14	20	231.46	7
org.jruby.ast	ArgsPushNode	2	7	7	0.67	7	1	1	26	642.14	9
org.jruby.ast.util	ArgsUtil	0	5	13	0	13	3	2.6	35	4082.78	5
org.jruby.ast	ArgumentNode	1	6	6	0.6	6	1	1	17	53.48	7
org.jruby.compiler	ArgumentsCallback	1	1	1	0	1	1	1	2	0	1
org.jruby.runtime	Arity	12	26	55	1.01	55	7	2.12	162	34258.83	37

entities

output

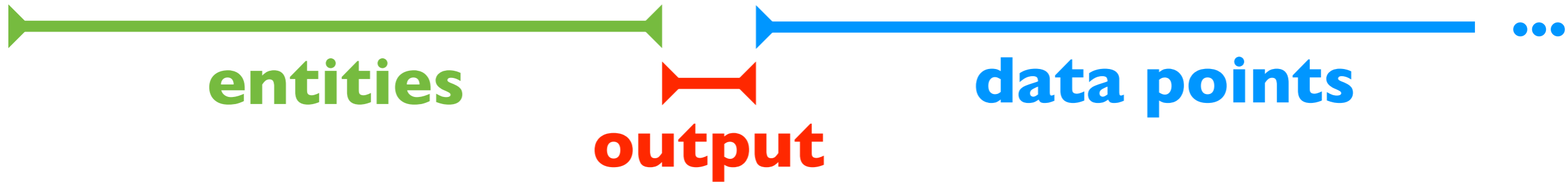
data points



Example Mining File

Package	Name	bugs	No. Methods	CCV	LCOM	TCC	MAXCC	AVCC	NOS	HEFF	UWCS
org.jruby.compiler	ASTCompiler	0	278	605	1	605	92	2.18	1574	424794.18	279
org.jruby.compiler	ASTCompiler\$OpElementAsgnArgumentsCallback	0	3	6	0	6	3	2	21	3592.03	4
org.jruby.compiler	ASTCompiler\$SpecificArityArguments	6	3	7	0.5	7	4	2.33	22	2660.26	5
org.jruby.compiler	ASTCompiler\$VariableArityArguments	0	3	3	0.5	3	1	1	7	31.61	4
org.jruby.compiler	ASTCompiler19	1	20	40	0	40	2	2	7	26264.39	20
org.jruby.compiler	ASTInspector	0	18	51	1.05	51	25	2.83	3	562739.87	45
org.jruby.evaluator	ASTInterpreter	4	22	55	0	55	6	2.5	9	38462.7	22
org.jruby.runtime	AbstractCompiledBlockCallback	0	0	0	0	0	0	0	0	0	0
org.jruby.ext.ffi	AbstractInvoker	5	7	11	0.92	11	4	1.57	3	6346.91	9
org.jruby.ext.ffi	AbstractMemory	0	33	105	0	105	5	1.54	31	87618.03	71
org.jruby.ast.executable	AbstractScript	0	166	187	0	187	3	1.13	41	28142.16	190
org.jruby.compiler.impl	AbstractVariableCompiler	8	20	61	0	61	12	2.75	27	129005.27	27
org.jruby.util	Adler32Ext	0	10	10	0.67	10	2	1.12	27	484	12
org.jruby.internal.runtime.methods	AliasMethod	0	10	45	0	45	1	1.12	24	1383.47	53
org.jruby.ast	AliasNode	0	7	7	0	7	1	1	3	90.68	9
org.jruby.ext.ffi	AllocatedNativeMemoryIO	2	2	2	0	2	1	1	3	0	2
org.jruby.ext.ffi.jffi	AllocatedNativeMemoryIO	3	5	8	0	8	2	1.6	20	1138.3	8
org.jruby.ext.ffi.jna	AllocatedNativeMemoryIO	0	2	2	0.67	2	2	1.5	19	884.91	5
org.jruby.ast	AndNode	0	8	8	0	8	2	1.14	27	871.77	9
org.jruby.anno	AnnotationBinder	0	3	3	1	3	1	1	27	32.73	5
org.jruby.anno	AnnotationBinders\$AnnotationBindingProcessor	0	14	14	0.88	14	6	2.8	41	2180.79	7
org.jruby.anno	AnnotationBinders\$AnnotationBinding	0	8	66	0.86	66	22	8.25	246	230539.24	10
org.jruby.ast	ArgAuxiliaryNode	0	6	6	0.8	6	1	1	17	85.47	8
org.jruby.ast	ArgsCatNode	2	7	7	0.5	7	1	1	26	806.1	9
org.jruby.ast	ArgsNode	10	2	2	0	2	1	1	8	0	2
org.jruby.ast	ArgsOneNode	4	34	64	0.9	64	7	1.88	146	33903.67	48
org.jruby.ast	ArgsOneNode	11	6	6	0	6	1	1	17	115.79	6
org.jruby.ast	ArgsOneNode	2	7	8	0	8	2	1.14	20	231.46	7
org.jruby.ast	ArgsPushNode	2	7	7	0.67	7	1	1	26	642.14	9
org.jruby.ast.util	ArgsUtil	0	5	13	0	13	3	2.6	35	4082.78	5
org.jruby.ast	ArgumentNode	1	6	6	0.6	6	1	1	17	53.48	7
org.jruby.compiler	ArgumentsCallback	1	1	1	0	1	1	1	2	0	1
org.jruby.runtime	Arity	12	26	55	1.01	55	7	2.12	162	34258.83	37

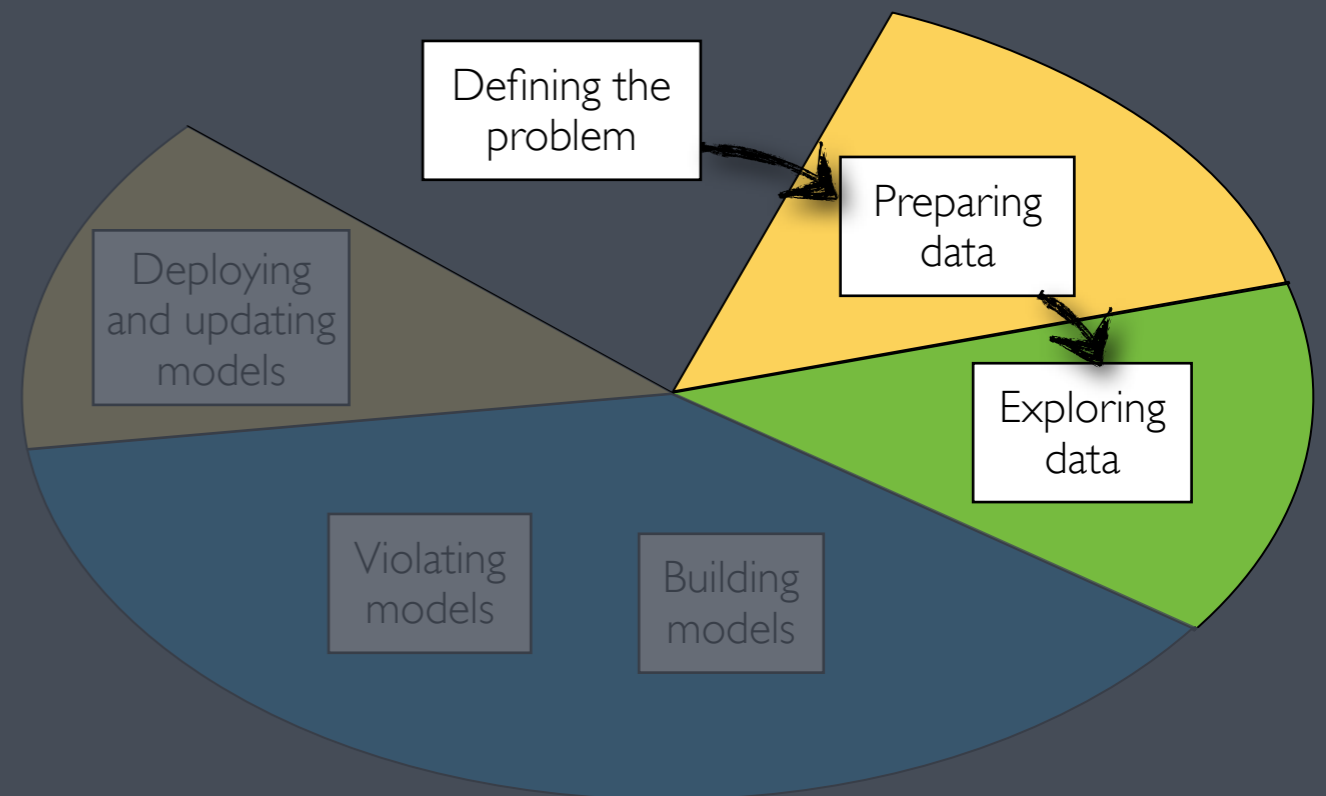
Careful! Large files!
e.g.: 5 million lines, 300 column



Step 3: Explore Data

You cannot validate the output if you don't know the input

- ▶ Descriptive data summary
 - ▶ max, min, mean, pareto, distribution
- ▶ Data Selection
 - ▶ Relevance of data
- ▶ Data reduction
 - ▶ aggregation, subset selection



Descriptive Data Summary

Step 1: Define the problem

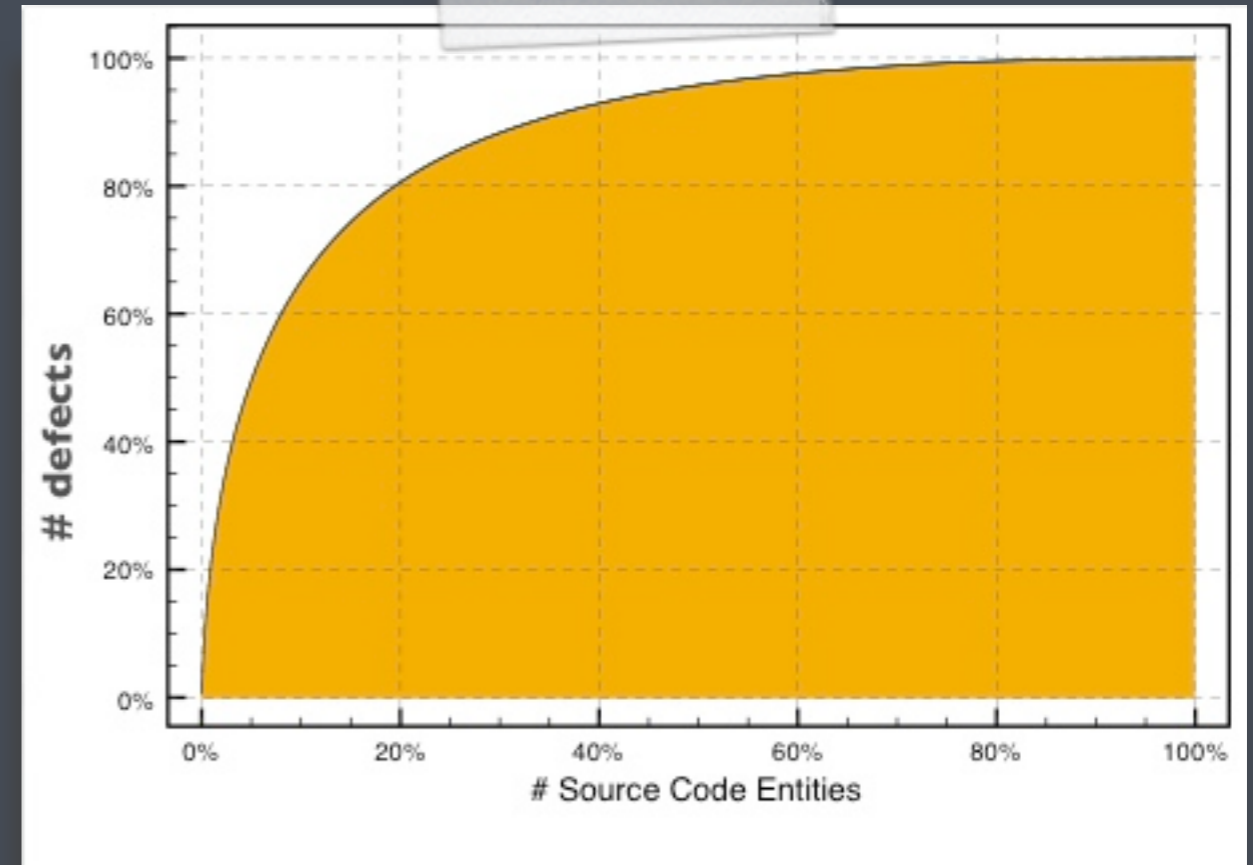
Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model

- ▶ How good can a prediction possibly be?
- ▶ Does it make sense to predict the top 20%

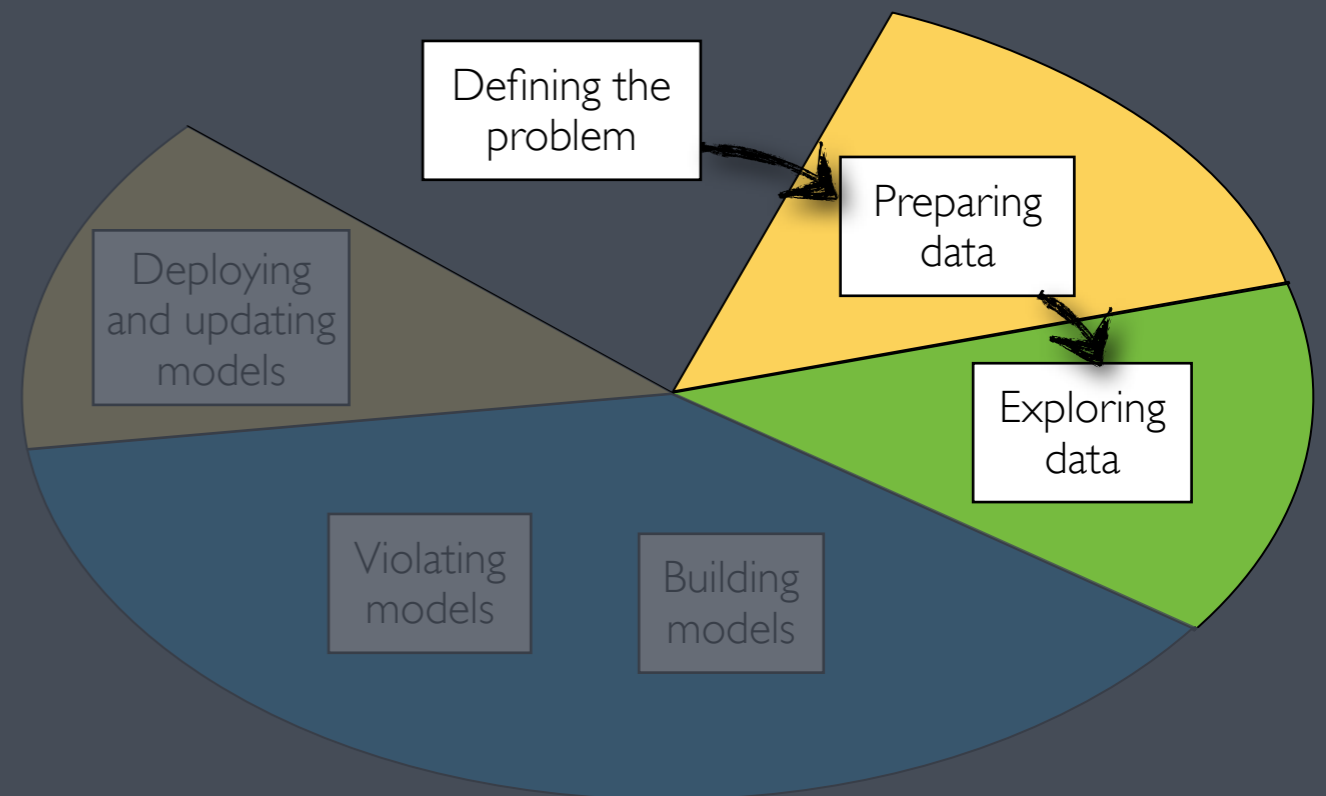


20% of entities contain 80% of defects

Step 3: Explore Data

Data sufficiency

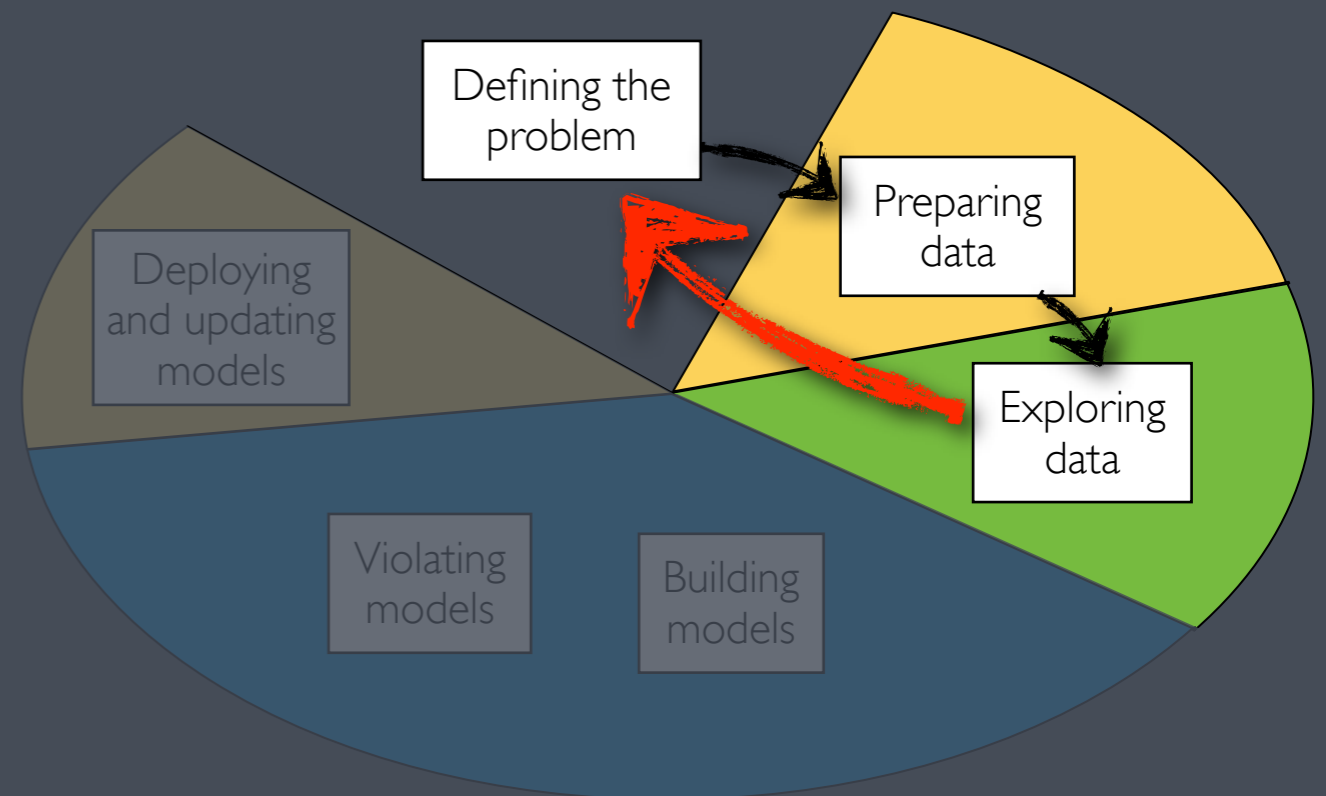
- ▶ Maybe the data will not help to solve the problem
- ▶ Redefine problem
- ▶ Search for alternatives
- ▶ Access different data



Step 3: Explore Data

Data sufficiency

- ▶ Maybe the data will not help to solve the problem
- ▶ Redefine problem
- ▶ Search for alternatives
- ▶ Access different data



Step 3: Explore Data

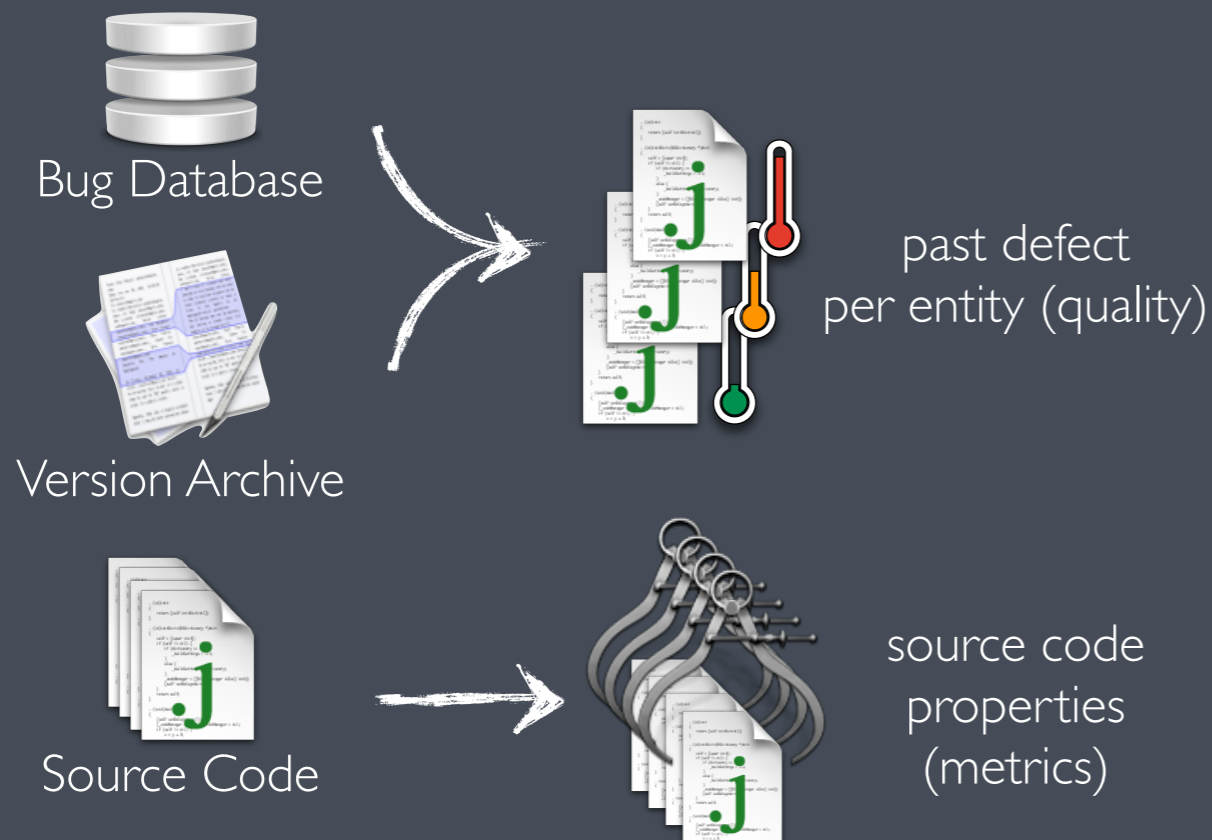
Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model



*Does complexity
correlate with defects?*

Step 3: Explore Data

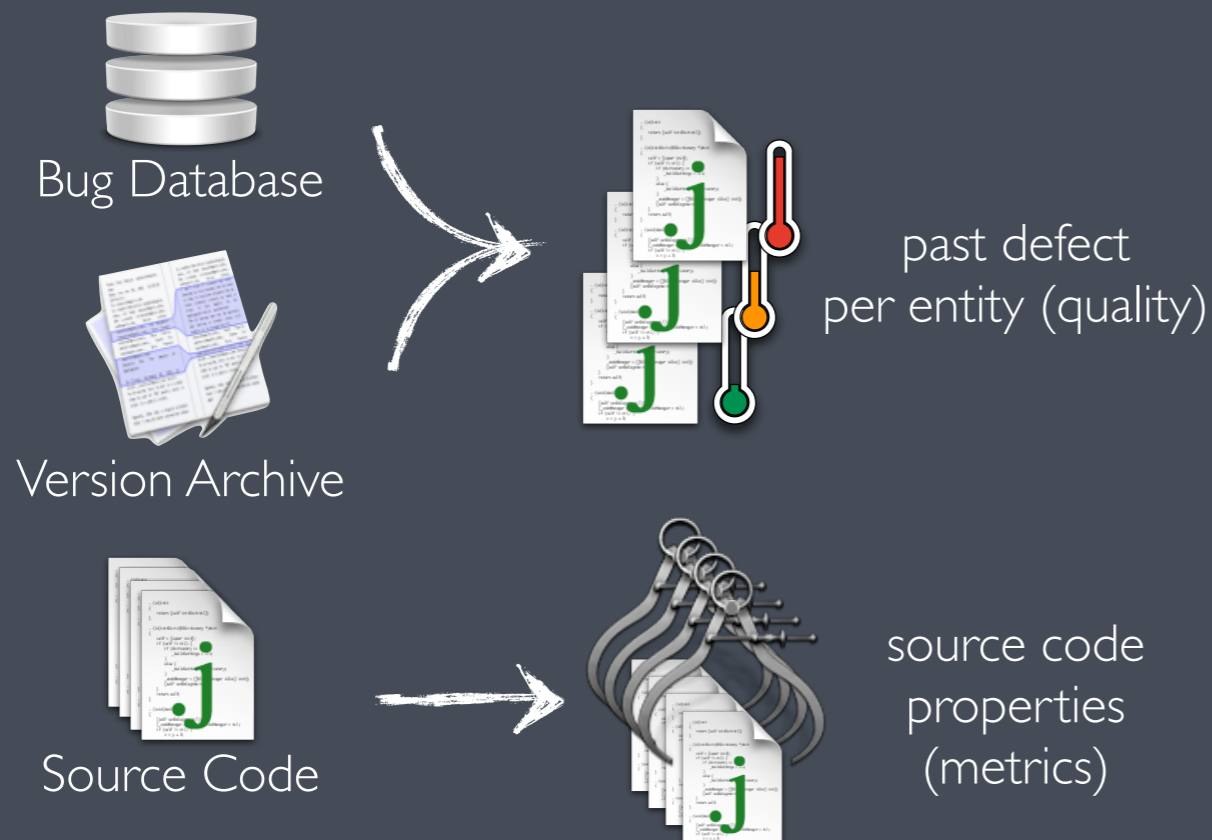
Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model

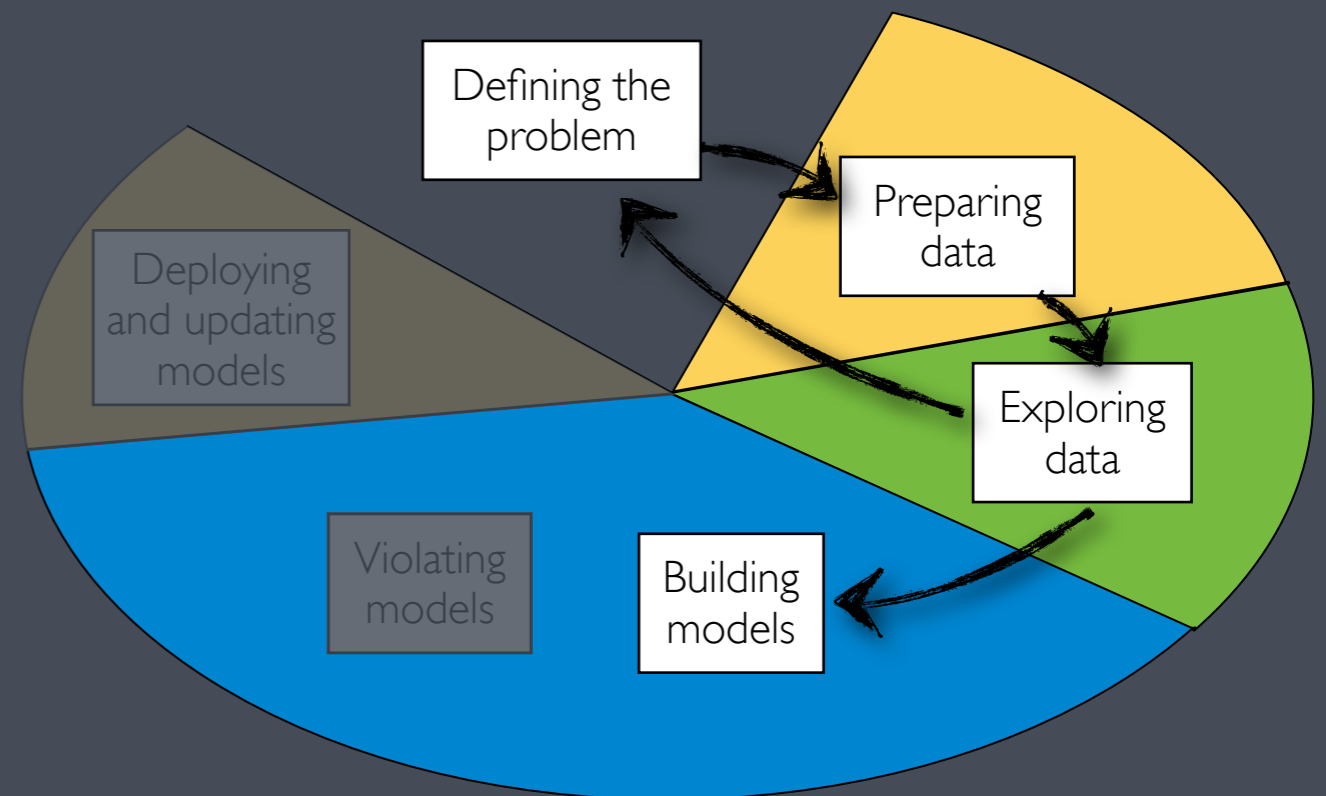


Does complexity correlate with defects?

YES!

Step 4: Build Model

- ▶ Mining model only container
 - ▶ parameters and mining structure
 - ▶ output value
- ▶ Now we need some statistics / machine learners



Example Mining File

Package	Name	bugs	No. Methods	CCV	LCOM	TCC	MAXCC	AVCC	NOS	HEFF	UWCS
org.jruby.compiler	ASTCompiler	0	278	605	1	605	92	2.18	1574	424794.18	279
org.jruby.compiler	ASTCompiler\$OpElementAsgnArgumentsCallback	0	3	6	0	6	3	2	21	3592.03	4
org.jruby.compiler	ASTCompiler\$SpecificArityArguments	6	3	7	0.5	7	4	2.33	22	2660.26	5
org.jruby.compiler	ASTCompiler\$VariableArityArguments	0	3	3	0.5	3	1	1	7	31.61	4
org.jruby.compiler	ASTCompiler19	1	20	40	0	40	10	2	117	26264.39	20
org.jruby.compiler	ASTInspector	0	18	51	1.05	51	25	2.83	433	562739.87	45
org.jruby.evaluator	ASTInterpreter	4	22	55	0	55	6	2.5	169	38462.7	22
org.jruby.runtime	AbstractCompiledBlockCallback	0	0	0	0	0	0	0	1	0	0
org.jruby.ext.ffi	AbstractInvoker	5	7	11	0.92	11	4	1.57	42	6346.91	9
org.jruby.ext.ffi	AbstractMemory	0	68	105	0.97	105	5	1.54	310	87618.03	71
org.jruby.ast.executable	AbstractScript	1	166	187	24	187	3	1.13	413	28142.16	190
org.jruby.compiler.impl	AbstractVariableCompiler	8	20	61	0.71	61	12	3.05	279	129005.27	27
org.jruby.util	Adler32Ext	0	9	10	0.67	10	2	1.11	27	484	12
org.jruby.internal.runtime.methods	AliasMethod	2	40	45	0.96	45	6	1.12	102	1383.47	53
org.jruby.ast	AliasNode	2	7	7	0.67	7	1	1	24	90.68	9
org.jruby.ext.ffi	AllocatedDirectMemoryIO	11	2	2	0	2	1	1	3	0	2
org.jruby.ext.ffi.jffi	AllocatedNativeMemoryIO	3	5	8	0.75	8	2	1.6	20	1138.3	8
org.jruby.ext.ffi.jna	AllocatedNativeMemoryIO	0	4	6	0.67	6	2	1.5	19	884.91	5
org.jruby.ast	AndNode	2	7	8	0.5	8	2	1.14	27	871.77	9
org.jruby.anno	AnnotationBinder	6	3	3	1	3	1	1	27	32.73	5
org.jruby.anno	AnnotationBinder\$AnnotationBindingProcessor	2	5	14	0.88	14	6	2.8	41	2180.79	7
org.jruby.anno	AnnotationBinder\$AnnotationBindingProcessor\$RubyClassVisitor	0	8	66	0.86	66	22	8.25	246	230539.24	10
org.jruby.ast	ArgAuxillaryNode	8	6	6	0.8	6	1	1	17	85.47	8
org.jruby.ast	ArgsCatNode	2	7	7	0.5	7	1	1	26	806.1	9
org.jruby.ast	ArgsNoArgNode	10	2	2	0	2	1	1	8	0	2
org.jruby.ast	ArgsNode	4	34	64	0.9	64	7	1.88	146	33903.67	48
org.jruby.ast	ArgsPreOneArgNode	11	6	6	0	6	1	1	17	115.79	6
org.jruby.ast	ArgsPreTwoArgNode	2	7	8	0	8	2	1.14	20	231.46	7
org.jruby.ast	ArgsPushNode	2	7	7	0.67	7	1	1	26	642.14	9
org.jruby.ast.util	ArgsUtil	0	5	13	0	13	3	2.6	35	4082.78	5
org.jruby.ast	ArgumentNode	1	6	6	0.6	6	1	1	17	53.48	7
org.jruby.compiler	ArgumentsCallback	1	1	1	0	1	1	1	2	0	1
org.jruby.runtime	Arity	12	26	55	1.01	55	7	2.12	162	34258.83	37

entities

output

data points



Building the Model

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model

▶ **Regression**

- ▶ Predicting concrete, continuous values
- ▶ Difficult and very imprecise
- ▶ But desirable

▶ **Classification**

- ▶ Predicting class labels (e.g. more than X defects or not)
- ▶ Easier and more precise
- ▶ Vague information (how many defects in code?)

Building the Model

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model

Building the Model

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model

Support Vector Machine

Rule-Based Classification

Linear Regression

Lazy Learners

Decision Tree

Bayesian Network

Logistic Regression



Training and Testing

▶ **Training set**

- ▶ The data set to train the model
- ▶ Which columns correlate with output values?
- ▶ Which columns correlate with each other?

▶ **Testing set**

- ▶ A data set independent of the training data set
- ▶ used to fine-tune the estimates of the model parameters



Training and Testing

Random split

- + **Only one version needed**
- + **No overlaps between training and testing entities**
- **Does not reflect real life**
- **Which random set is the best one?** (because they are all different)

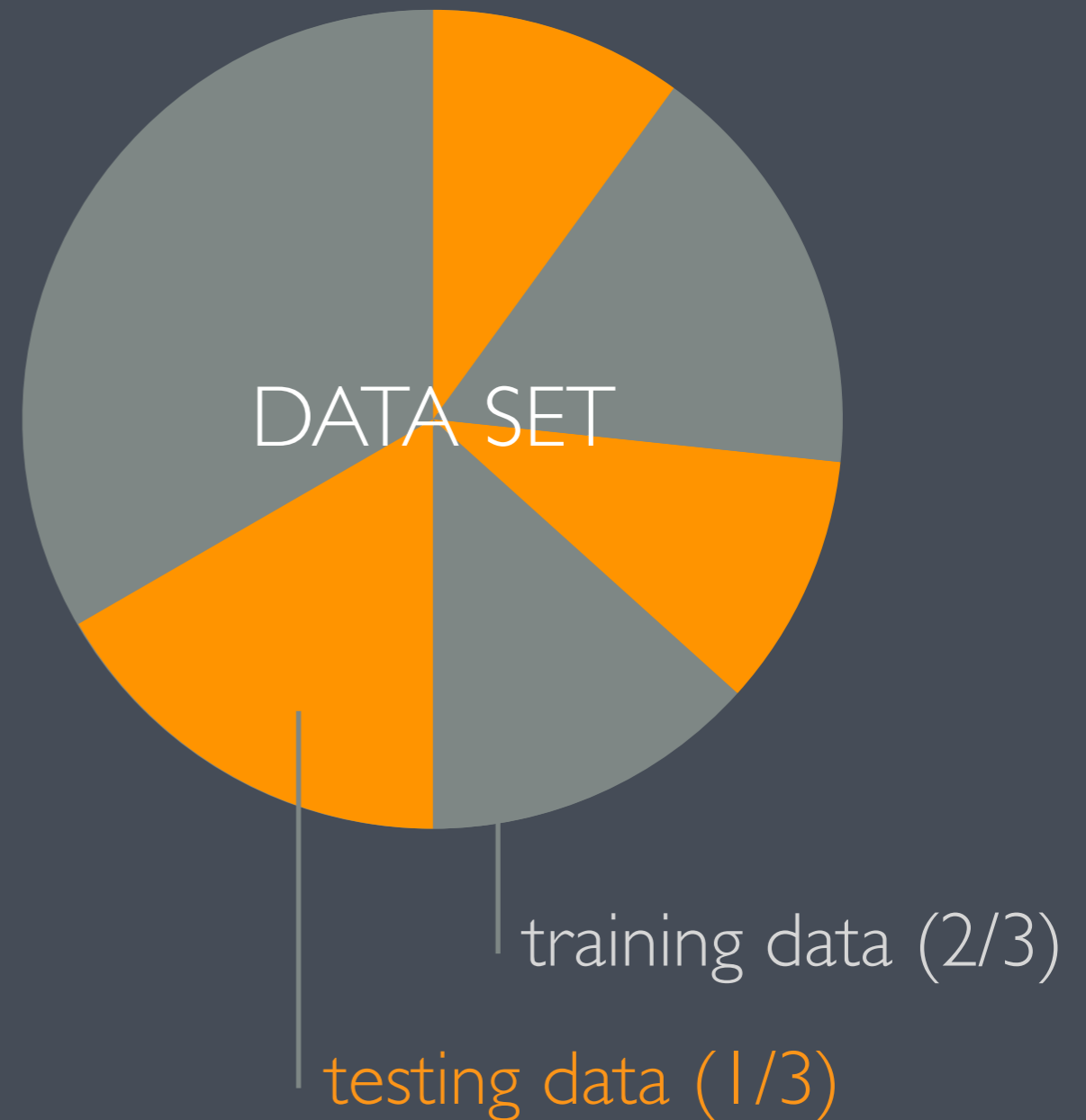


DATA SET

Training and Testing

Random split

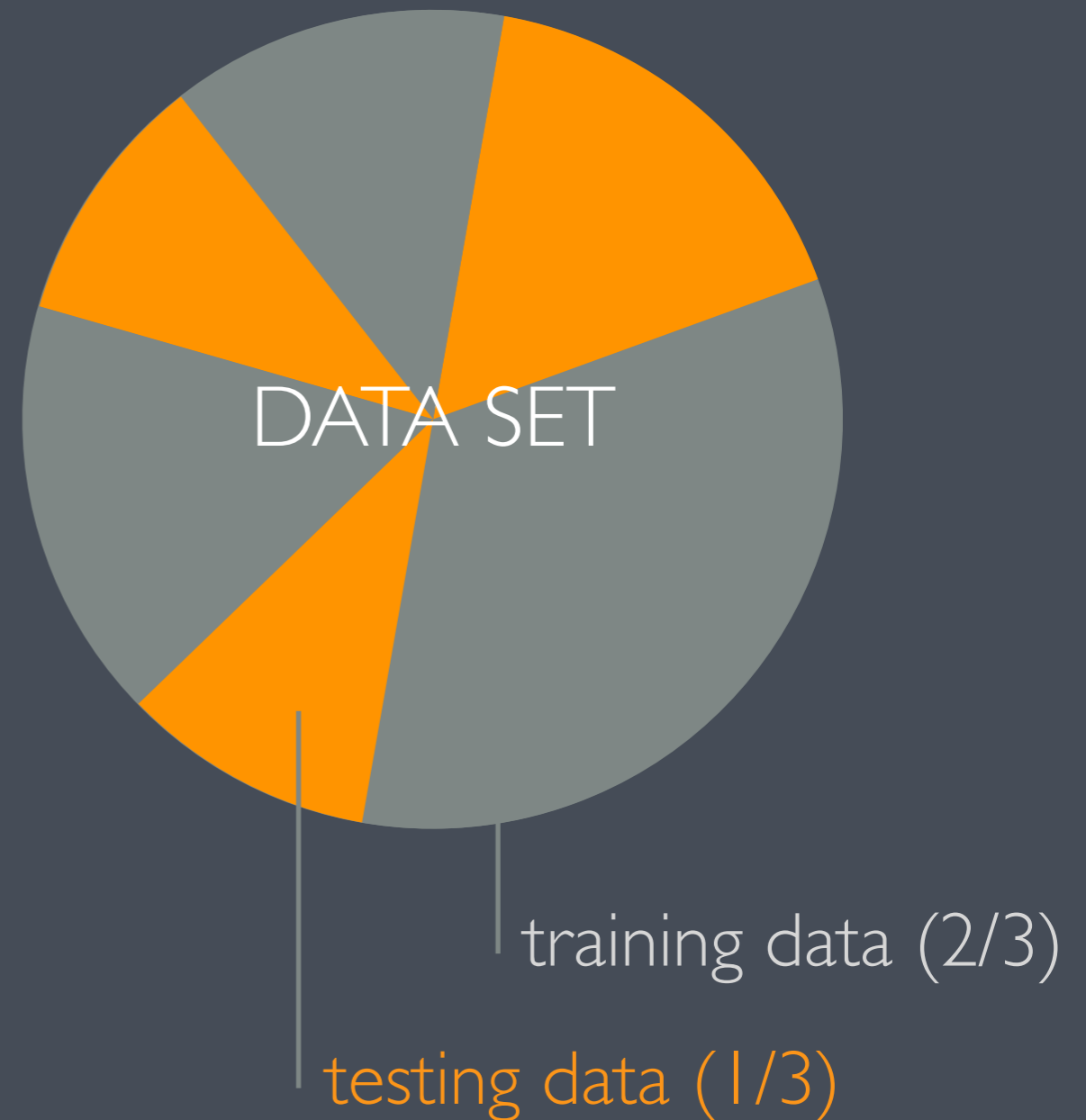
- + **Only one version needed**
- + **No overlaps between training and testing entities**
- **Does not reflect real life**
- **Which random set is the best one? (because they are all different)**



Training and Testing

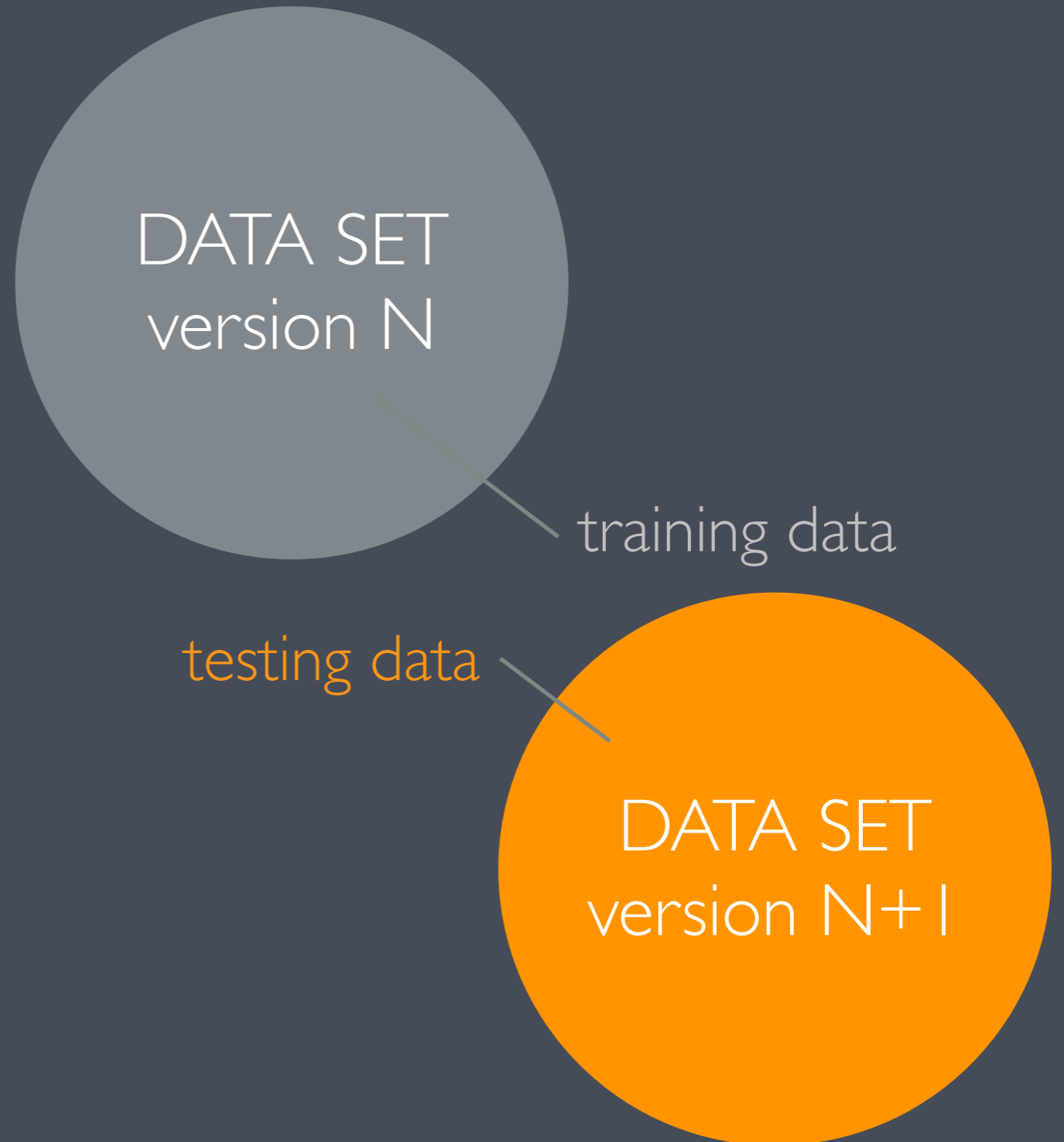
Random split

- + **Only one version needed**
- + **No overlaps between training and testing entities**
- **Does not reflect real life**
- **Which random set is the best one? (because they are all different)**



Training and Testing

- Forward estimation**
- + **Reflects real life**
- + **Reproducible result**
- **Two versions needed**



Step 4: Build Model

Step 4: Build Model



training set

Step 4: Build Model

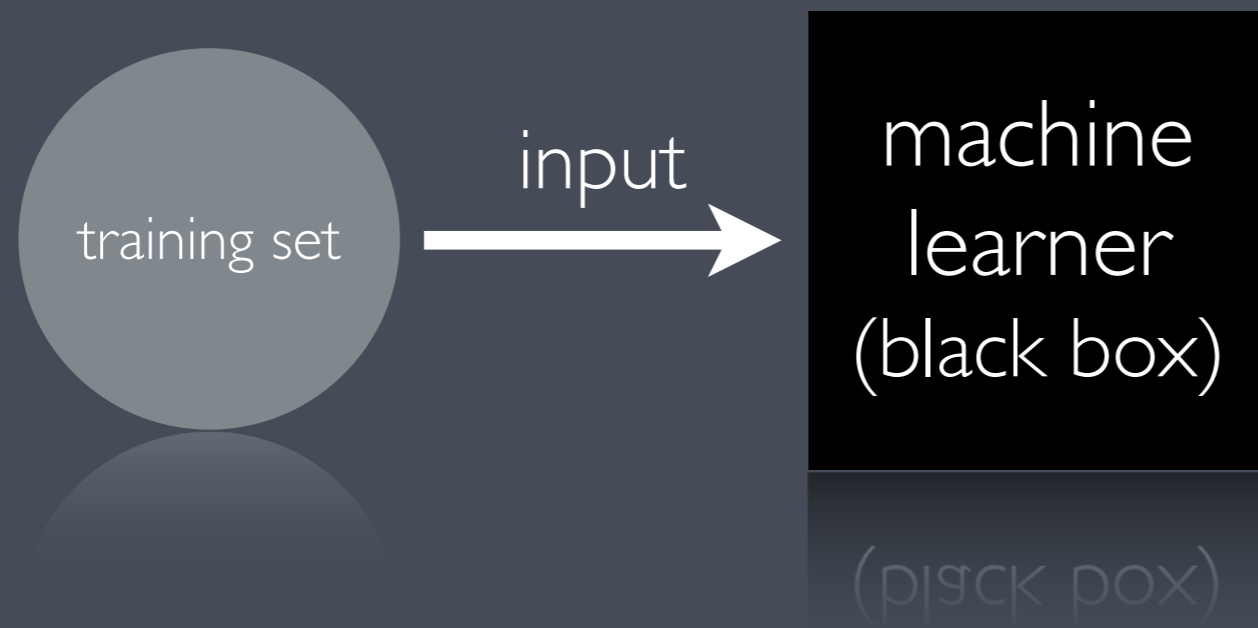
training set

machine
learner
(black box)

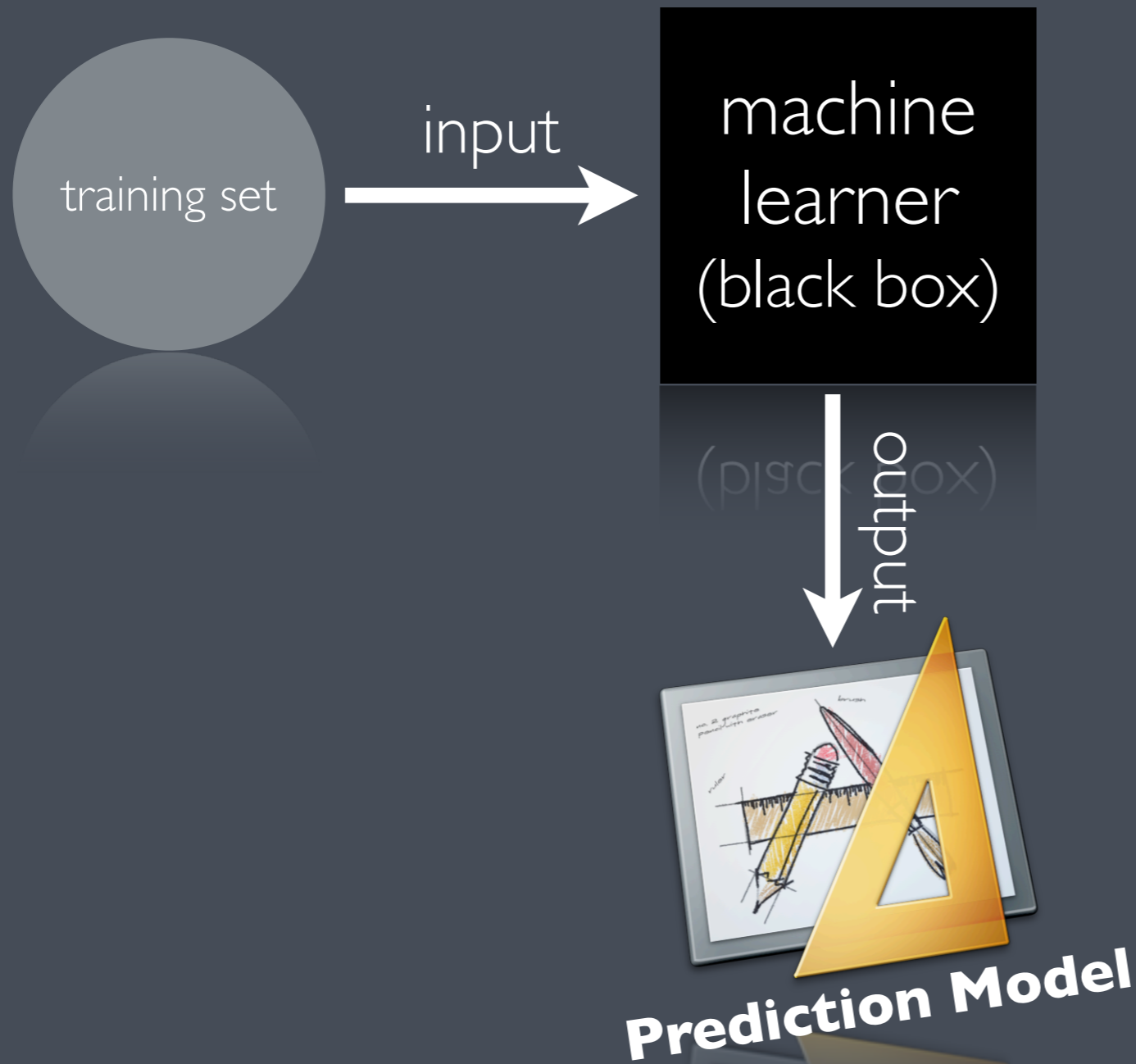
(black box)



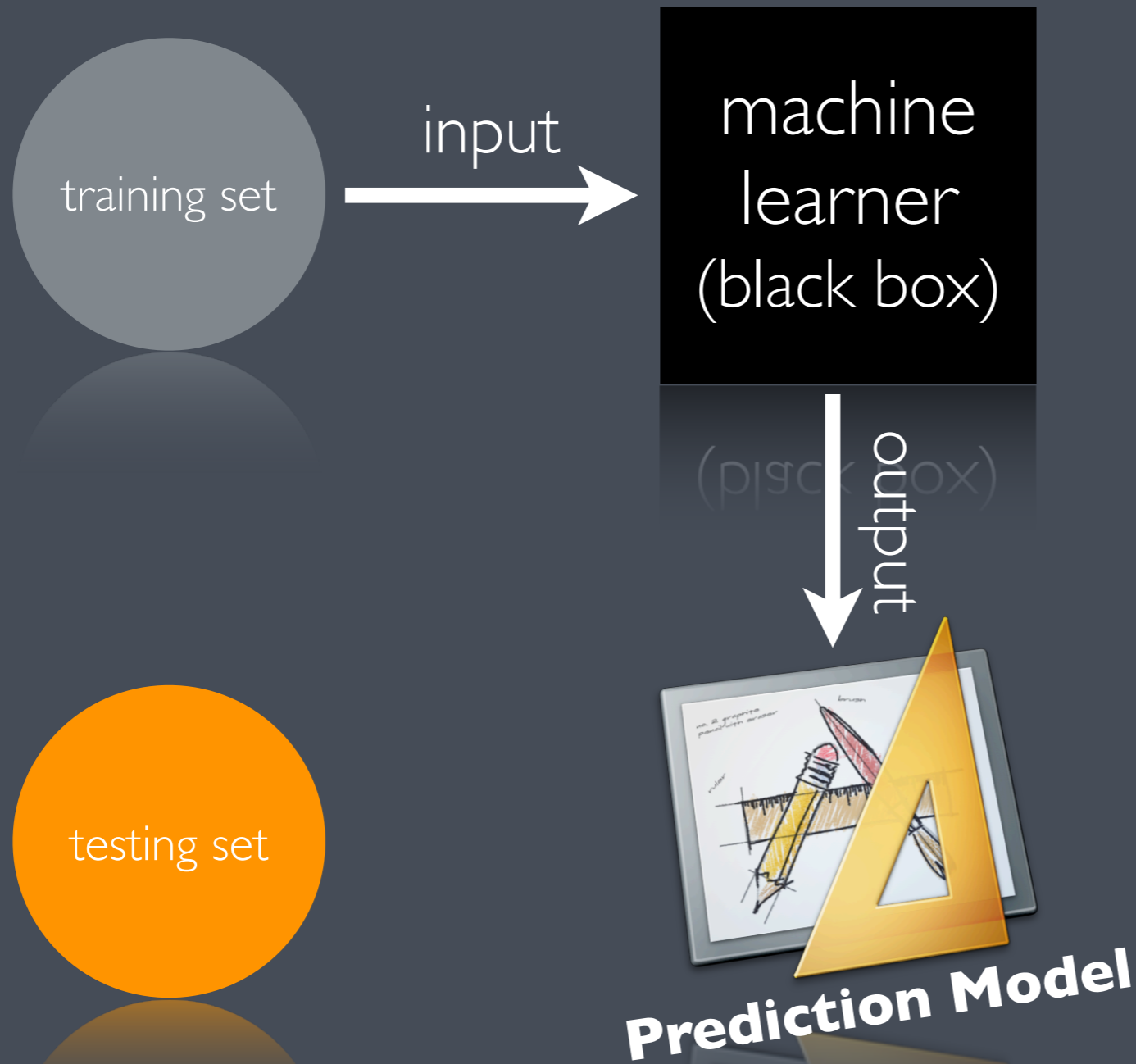
Step 4: Build Model



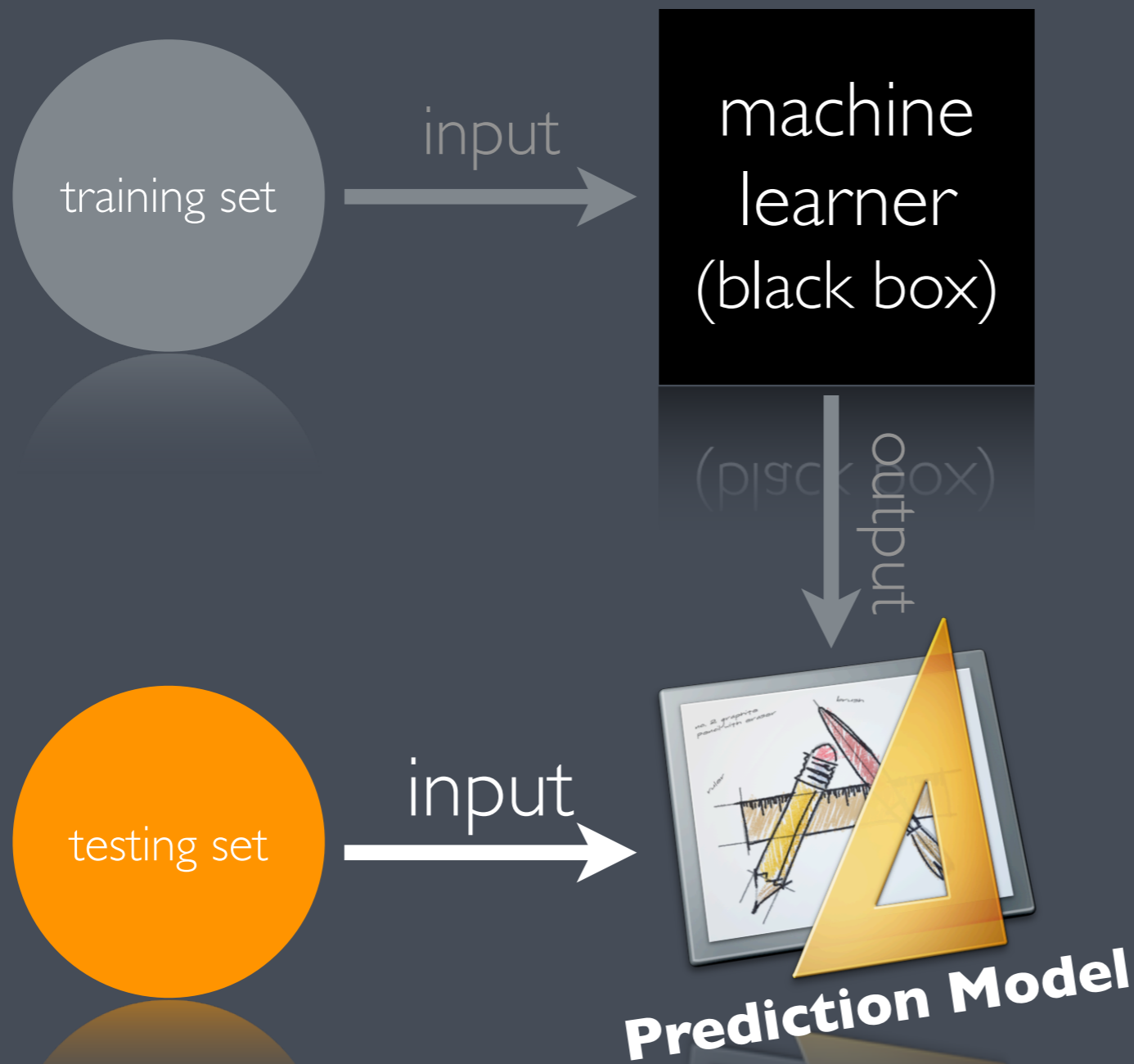
Step 4: Build Model



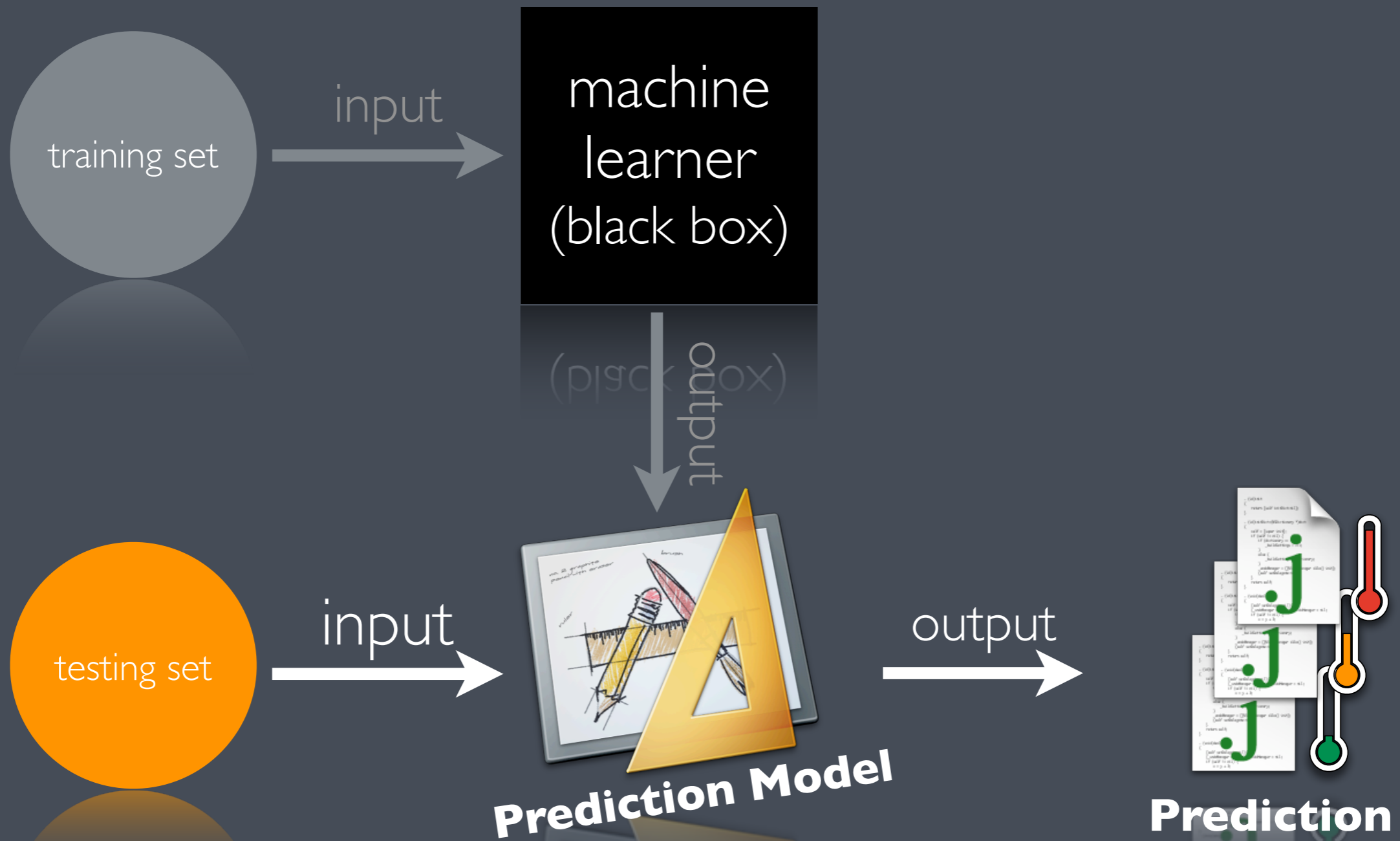
Step 4: Build Model



Step 4: Build Model

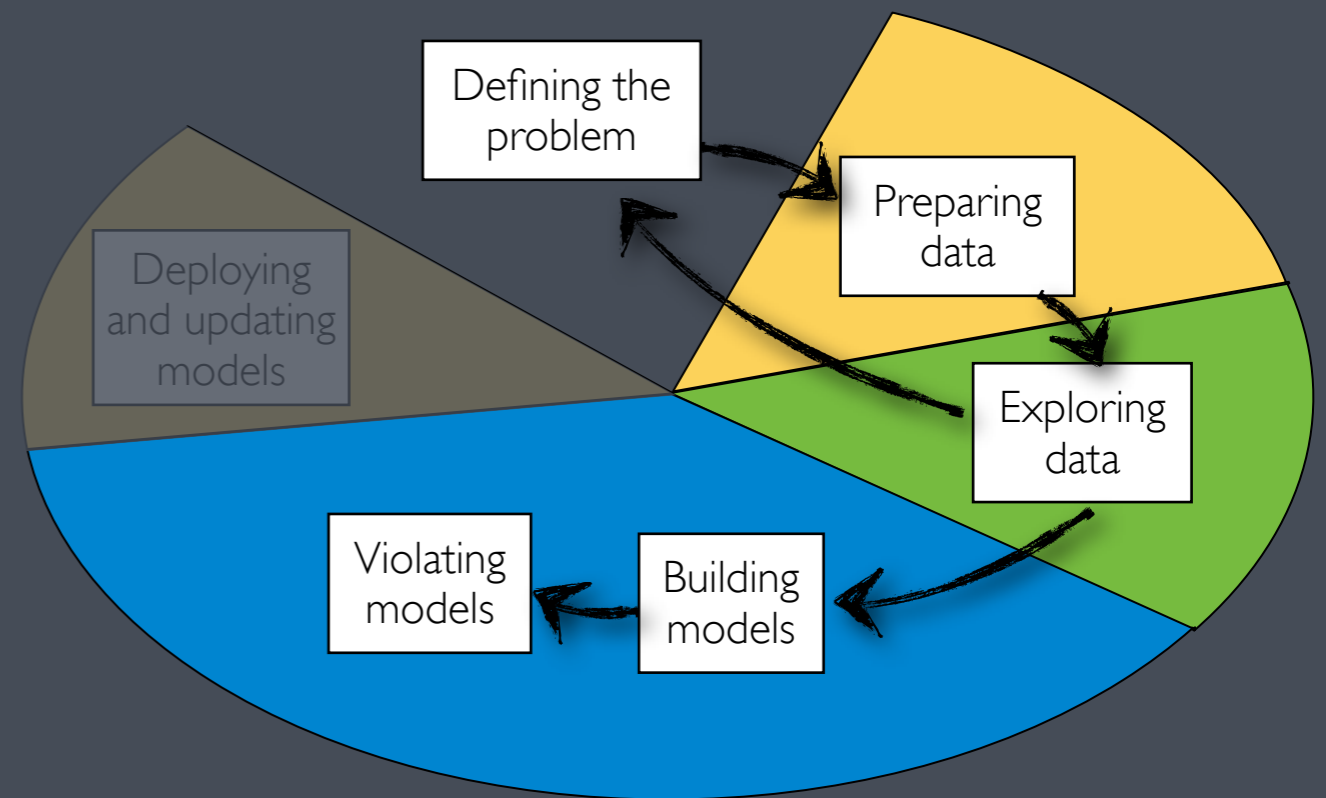


Step 4: Build Model



Step 5: Validating Model

- ▶ Test data has same structure but different content
- ▶ Goal is to use model to correctly estimate output values
- ▶ Compare estimation with real values (fine tuning)



Evaluation

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model

Evaluation

Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model



Never predict concrete number!

Because people will take them for real!

Evaluation

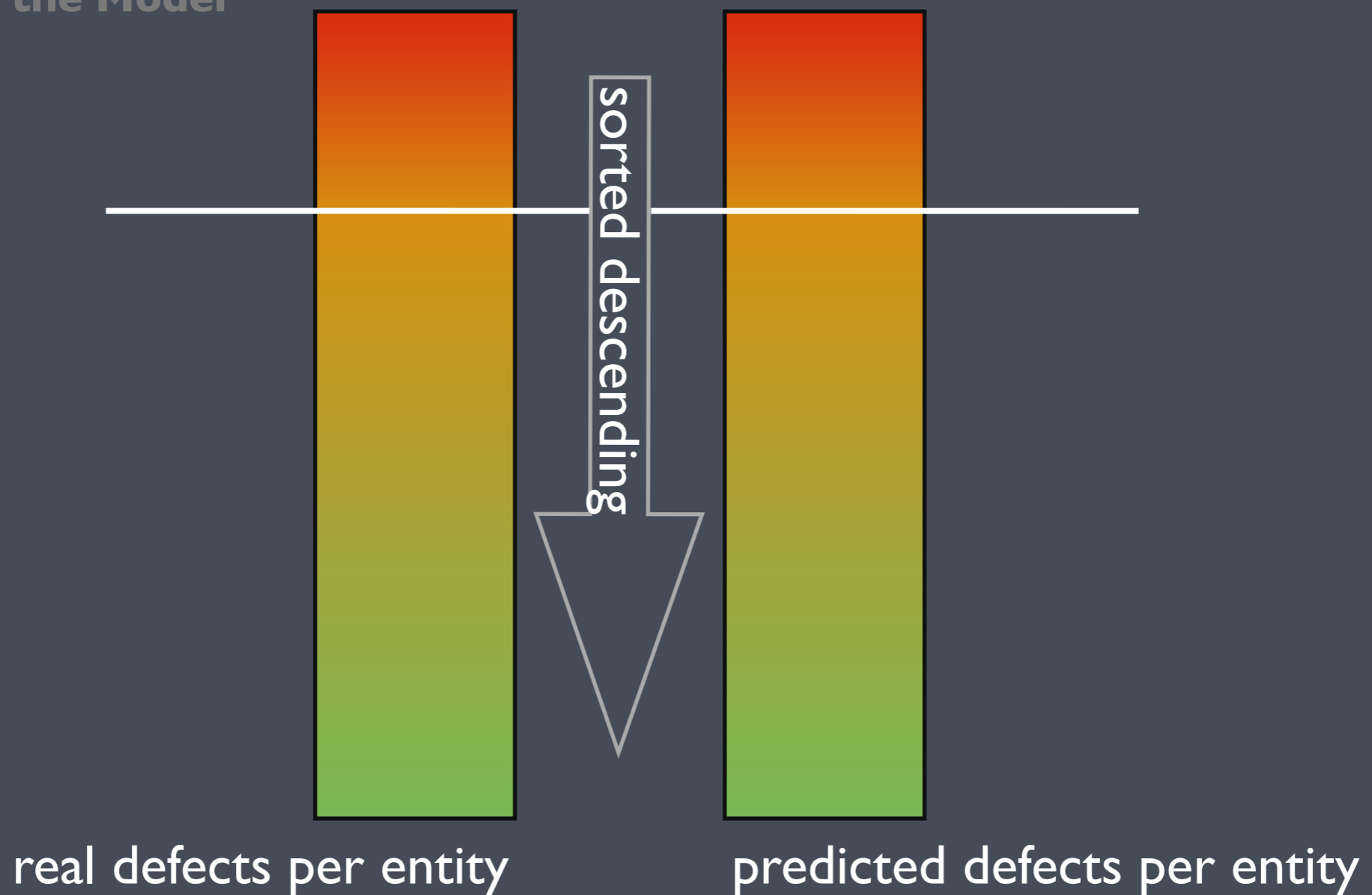
Step 1: Define the problem

Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model



Evaluation

Step 1: Define the problem

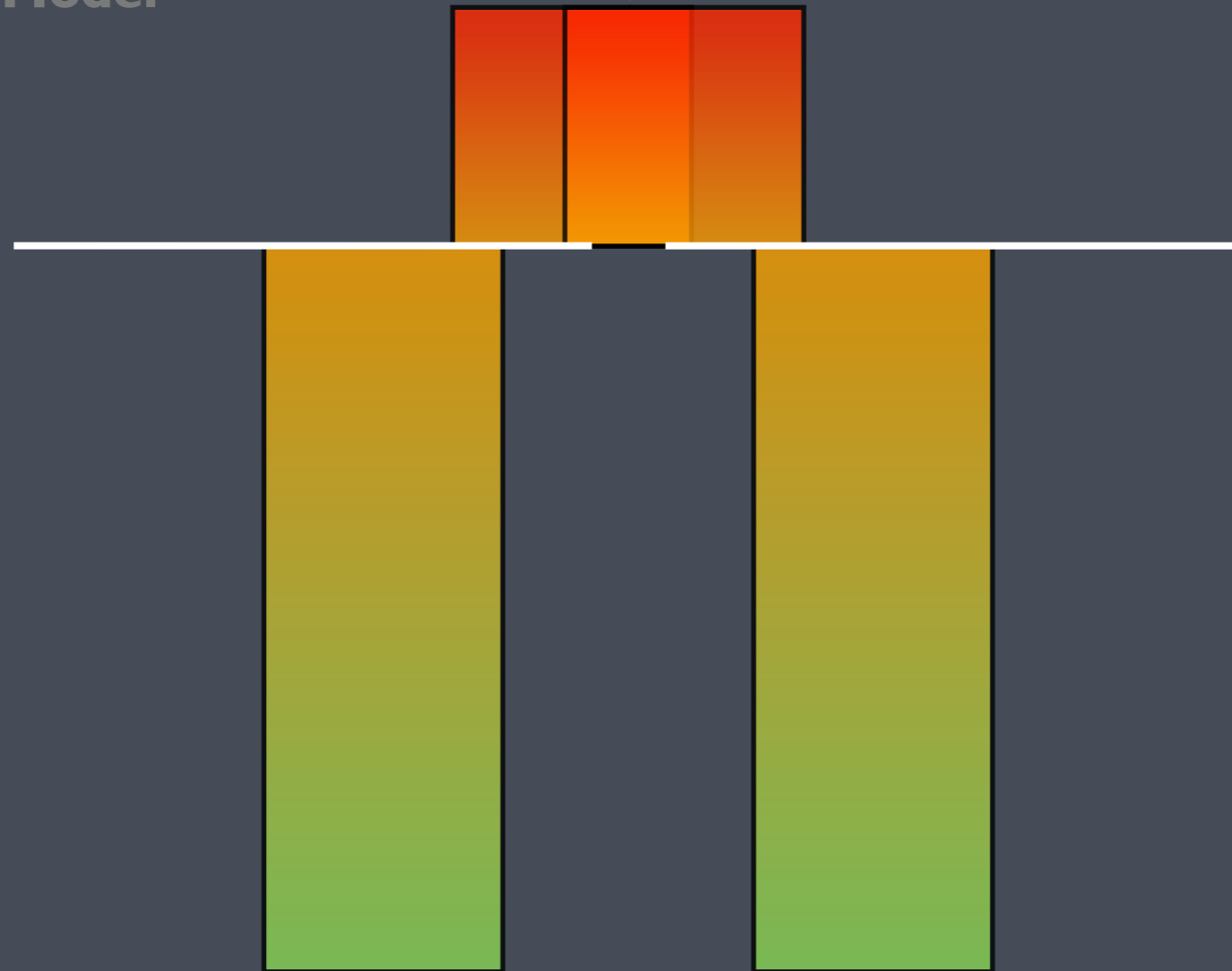
Step 2: Prepare Data

Step 3: Explore Data

Step 4: Building the Model

Step 5: Validating the Model

correctly predicted defect prone modules
(true positives)



real defects per entity

predicted defects per entity

Recall, Precision, Accuracy

		Predict defects ?	
		Yes	No
Real defects?	Yes	true positives	false negatives
	No	false positives	true negatives



Recall, Precision, Accuracy

		Predict defects ?	
		Yes	No
Real defects?	Yes	true positives	false negatives
	No	false positives	true negatives



Precision

$$\frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

*Predicted defect prone entities
will be defect prone!*



Recall, Precision, Accuracy

		Predict defects ?	
		Yes	No
Real defects?	Yes	true positives	false negatives
	No	false positives	true negatives



Recall

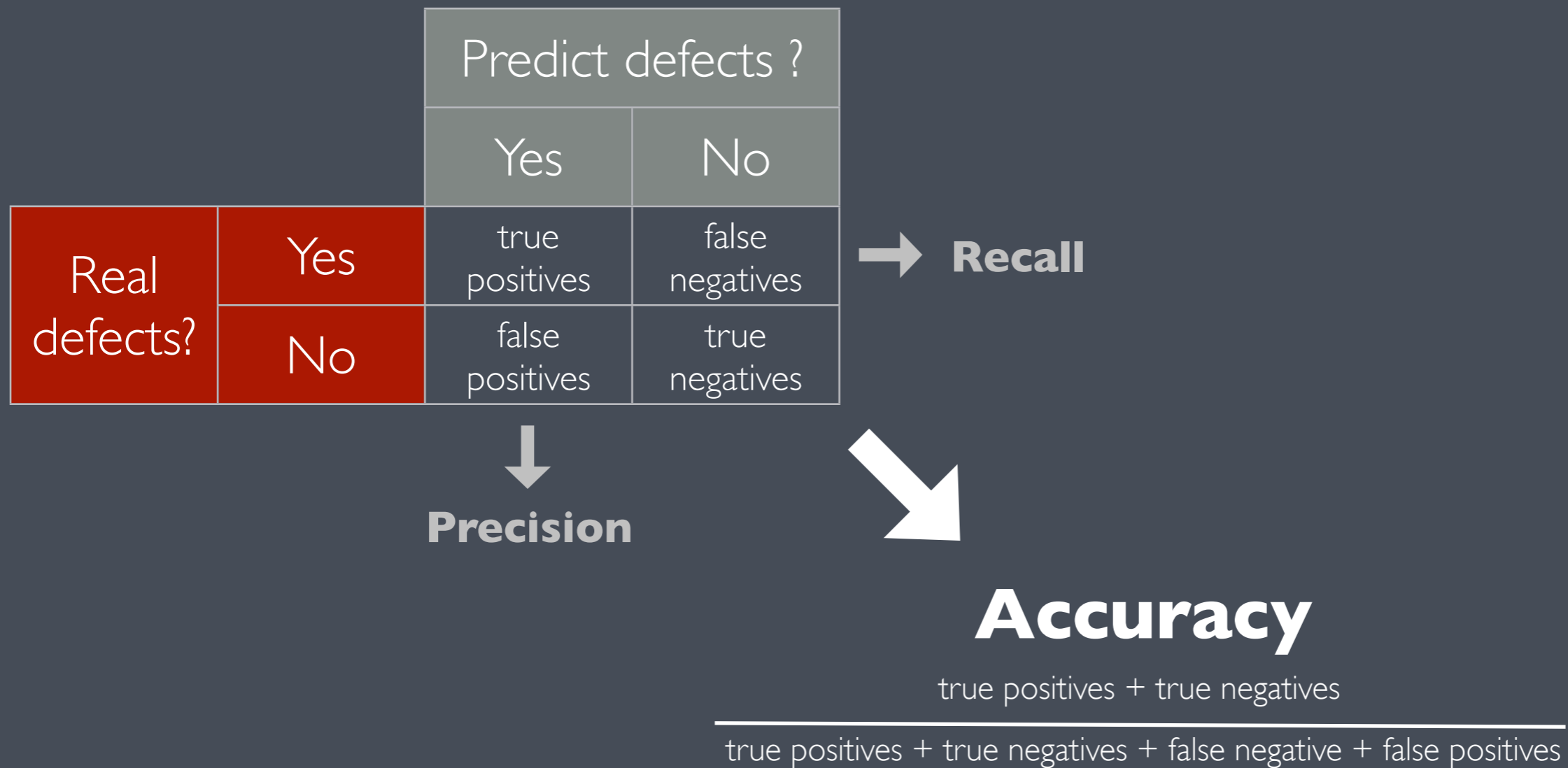
↓
Precision

$$\frac{\text{true positives}}{\text{true positives} + \text{false negative}}$$

*All defect prone entities
get predicted as defect prone.*



Recall, Precision, Accuracy

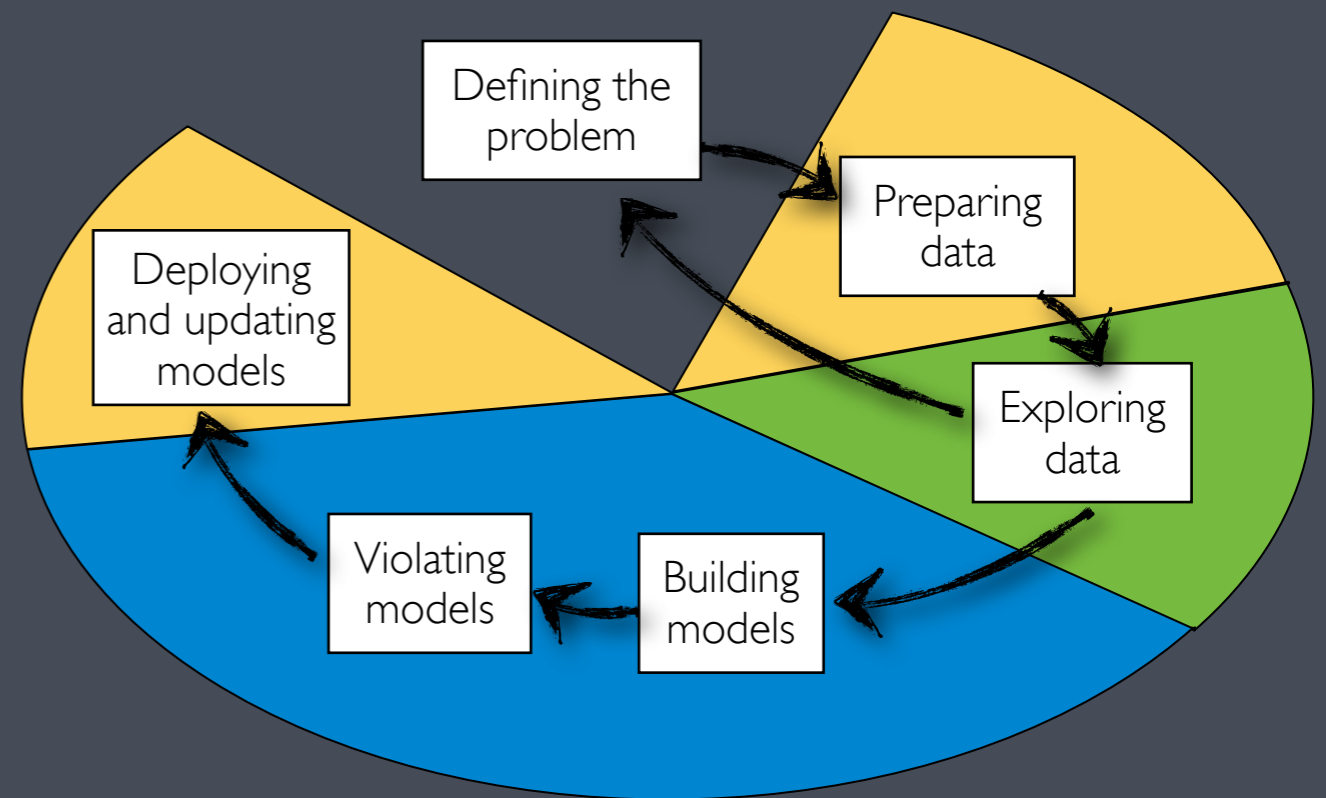


The overall correctness



Step 6: Deploying Model

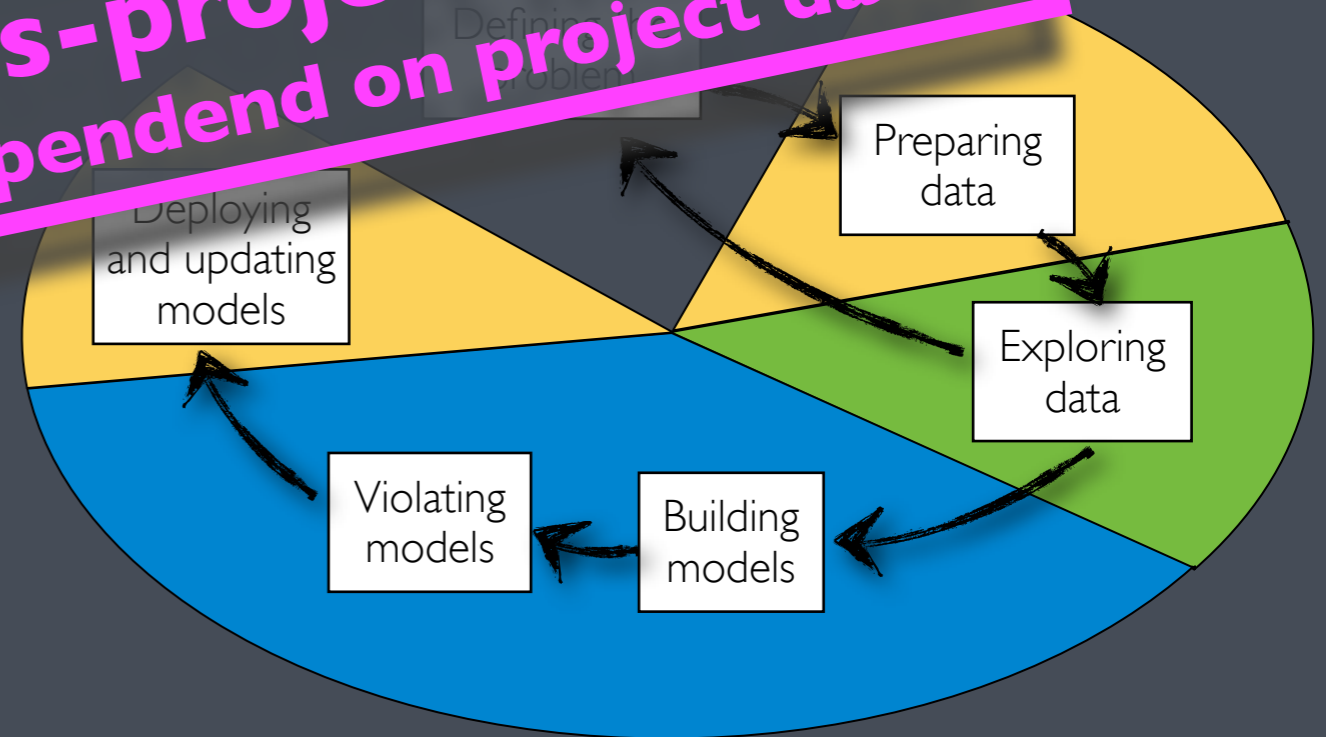
- ▶ Integrate model into development or quality assurance process
- ▶ Update model frequently (because change happens)
- ▶ Frequently validate the precision of your model



Step 6: Deploying Model

- ▶ Integrate model into development or quality assurance process
- ▶ Update model frequently (because change happens)
- ▶ Frequently validate the precision of your model

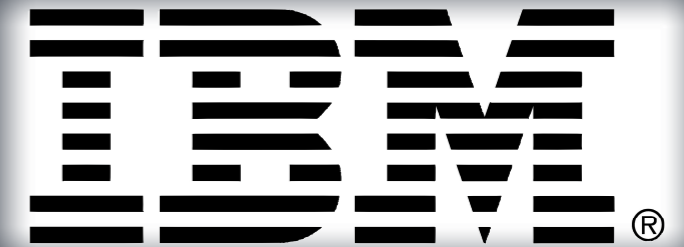
**Careful with cross-project models!
Many models highly dependend on project data!**



State of the Art



State of the Art



Prediction Results

Training	Testing	Precision	Recall	Accuracy
2.0	2.0	0.692	0.265	0.876
	2.1	0.478	0.191	0.890
	3.0	0.613	0.171	0.861
2.1	2.0	0.664	0.203	0.870
	2.1	0.668	0.160	0.900
	3.0	0.717	0.139	0.864
3.0	2.0	0.578	0.277	0.866
	2.1	0.528	0.220	0.894
	3.0	0.675	0.224	0.869



Predicting Defects for Eclipse [Revised for Dataset Version 2.0a]

Thomas Zimmermann
Saarland University
tz@acm.org

Rahul Premraj
Saarland University
premrj@cs.uni-sb.de

Andreas Zeller
Saarland University
zeller@acm.org

Predicting java classes: Classification:
has bugs, has no bugs



Prediction Results

Training	Testing	Precision	Recall	Accuracy
2.0	2.0	0.692	0.265	0.876
	2.1	0.478	0.191	0.890
	3.0	0.613	0.171	0.861
2.1	2.0	0.664	0.203	0.870
	2.1	0.668	0.160	0.900
	3.0	0.717	0.139	0.864
3.0	2.0	0.578	0.277	0.866
	2.1	0.528	0.220	0.894
	3.0	0.675	0.224	0.869



Predicting Defects for Eclipse [Revised for Dataset Version 2.0a]

Thomas Zimmermann
Saarland University
tz@acm.org

Rahul Premraj
Saarland University
premrj@cs.uni-sb.de

Andreas Zeller
Saarland University
zeller@acm.org

Predicting java classes: Classification:
has bugs, has no bugs

Prediction Results

Training	Testing	Precision	Recall	Accuracy
2.0	2.0	0.692	0.265	0.876
	2.1	0.478	0.191	0.890
	3.0	0.613	0.171	0.861
2.1	2.0	0.664	0.253	0.870
	2.1	0.688	0.160	0.900
	3.0	0.717	0.139	0.864
3.0	2.0	0.578	0.277	0.866
	2.1	0.528	0.220	0.894
	3.0	0.675	0.224	0.869

Complexity causes defects!



Predicting Defects for Eclipse [Revised for Dataset Version 2.0a]

Thomas Zimmermann
Saarland University
tz@acm.org

Rahul Premraj
Saarland University
premraj@cs.uni-sb.de

Andreas Zeller
Saarland University
zeller@acm.org

Predicting java classes: Classification:
has bugs, has no bugs

Prediction Results

Training	Testing	Precision	Recall	Accuracy
	2.0	0.692	0.265	0.876
2.0	2.1	0.478	0.191	0.890
	3.0	0.613	0.171	0.861
2.1	2.0	0.664	0.253	0.879
	2.1	0.717	0.280	0.884
3.0	2.0	0.528	0.277	0.866
	3.0	0.675	0.224	0.869

Complexity causes defects!

But not all defects come from complexity!



Predicting Defects for Eclipse
[Revised for Dataset Version 2.0a]

Thomas Zimmermann
Saarland University
tz@acm.org

Rahul Premraj
Saarland University
premraj@cs.uni-sb.de

Andreas Zeller
Saarland University
zeller@acm.org

Predicting java classes: Classification:
has bugs, has no bugs

What to mine?

Profiles

Bug Reports

e-mail

Code

Changes

What to mine?

Traces

Effort

Specification

Chats

Tests

Navigation

Models

Code

e-mail

Bug Reports

Changes

Profiles

Traces

Effort

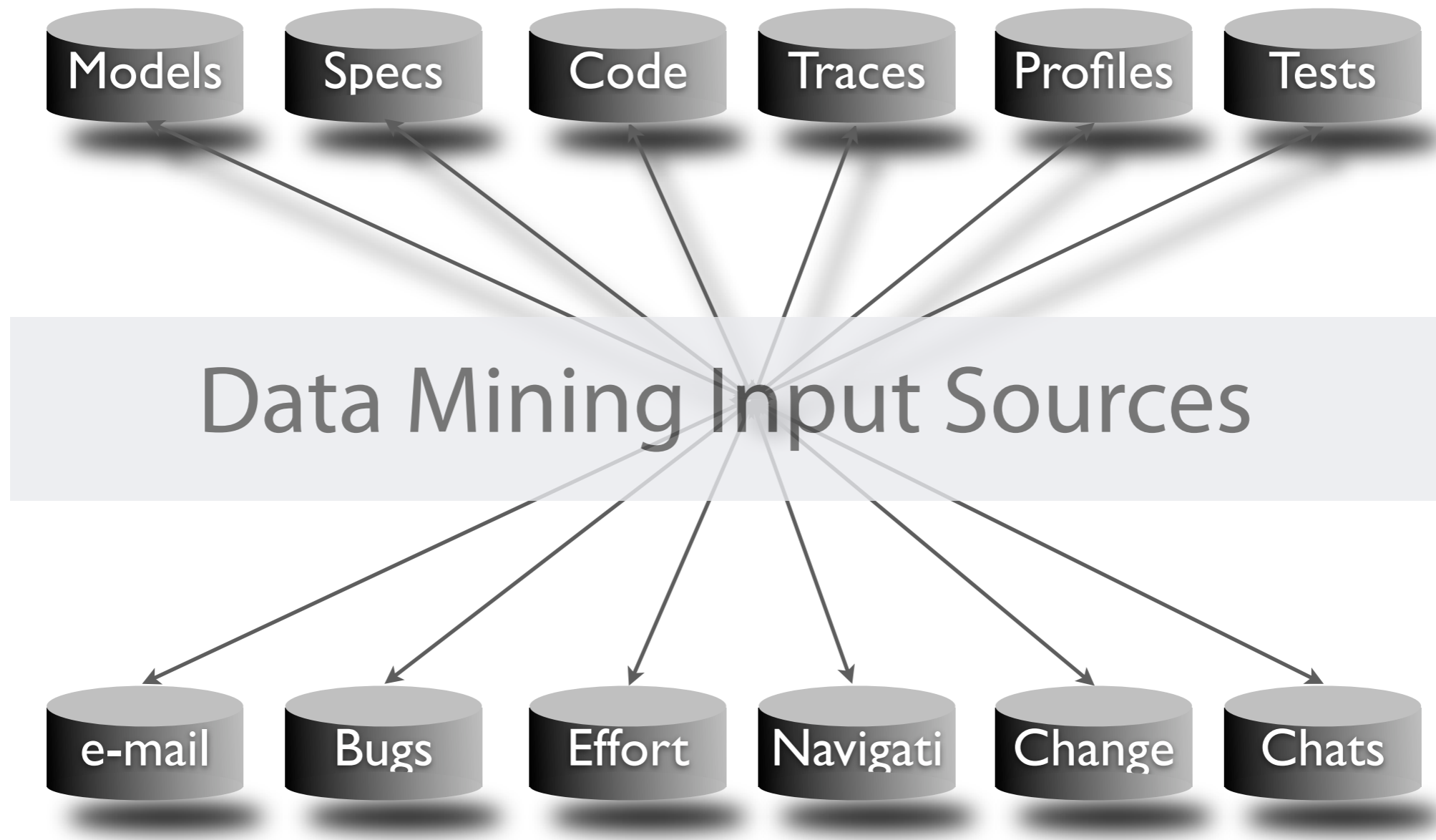
Chats

Sepcification

Tests

Navigation

Models





**People who changed function
f() also changed**





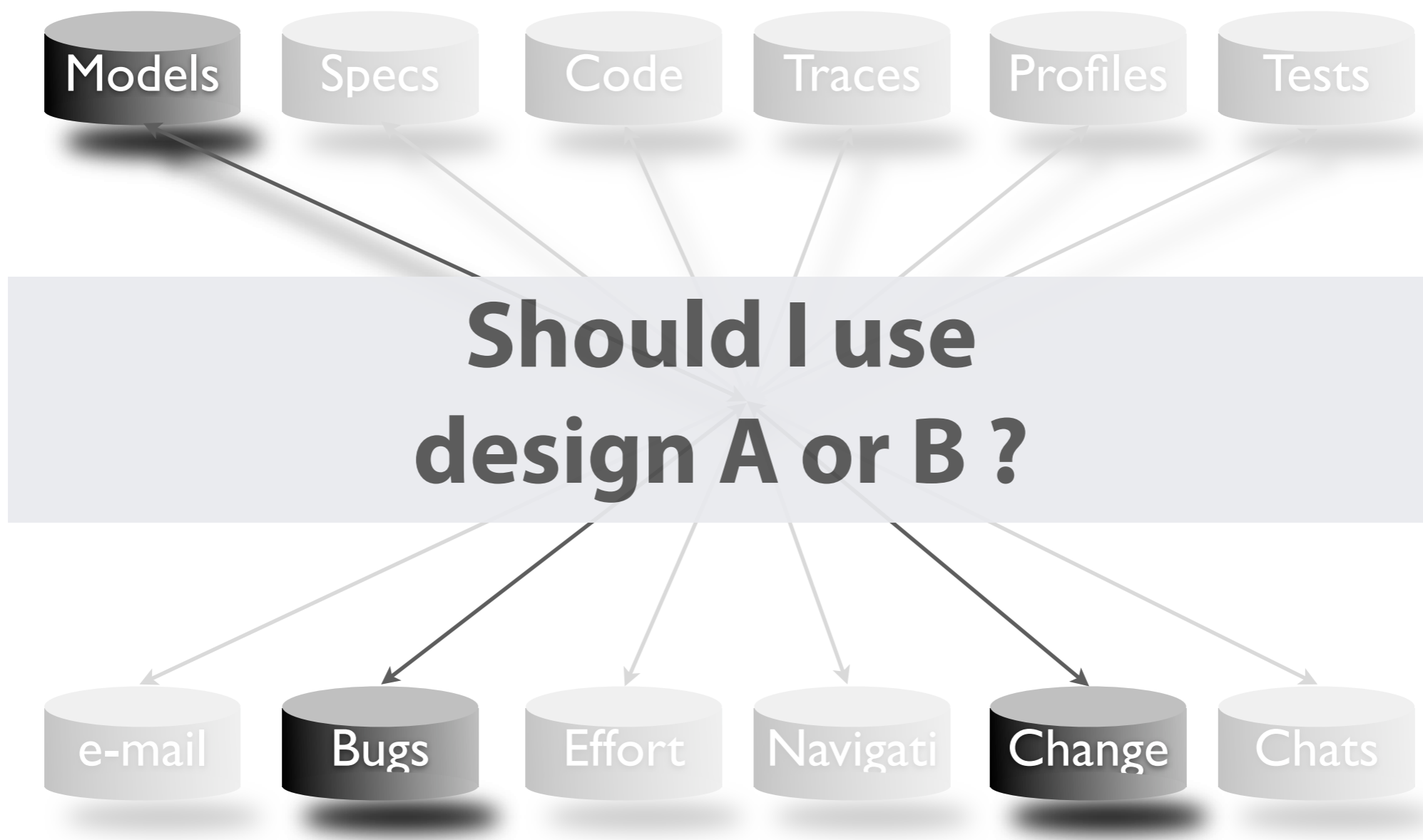
**Which modules should
I test most?**





**How long will it take
to fix this bug?**







**This requirement
is risky!**



Assistance

The screenshot displays the Team Concert IDE interface. The main window shows a defect view for 'Defect 61'. The defect details are as follows:

- Summary:** Test test.project.tests.TestProjectTest.testDoubles failed
- Type:** Defect
- Severity:** Major
- Found In:** Unassigned
- Created:** Feb 22, 2008 11:02 AM
- Created By:** Kim Herzig
- Team Area:** Bugs_and_Links Team / Bugs_and_Links
- Category:** Bugs_and_Links
- Tags:**
- Owned By:** Kim Herzig
- Priority:** High
- Planned For:** M1
- Estimate:**
- Due:** None
- Resolved:** Feb 22, 2008 11:03 AM
- Resolved By:** Kim Herzig

The description contains a stack trace for a `java.lang.AssertionError` with the message: `expected:<2.0> but was:<1.0>`. The stack trace includes the following frames:

```
at junit.framework.Assert.fail(Assert.java:47)
at junit.framework.Assert.failNotEquals(Assert.java:280)
at junit.framework.Assert.assertEquals(Assert.java:64)
at junit.framework.Assert.assertEquals(Assert.java:71)
...
```

The discussion section shows a comment from Kim Herzig, dated Feb 22, 2008, 11:02 AM, which repeats the stack trace.

The sidebar on the left shows work item statistics and an event log. The event log contains the following entries:

- Team Area created: Collaboration_Metric
- Kim Herzig created stream 'Collaboratio
- Project Area created: Collaboration_Metr

The bottom of the interface shows a table with the following columns: Id, Status, P, S, Summary, Owned By, Created By. The table is currently empty.

Assistance

Future environments will

- mine patterns from program + process
- apply rules to make predictions
- provide assistance in all development decisions
- adapt advice to project history



Empirical SE 2.0

Wikis

Joy of Use

Participation

Usability

Recommendation

Social Software

Collaboration

Perpetual Beta

Simplicity

Empirical SE 2.0

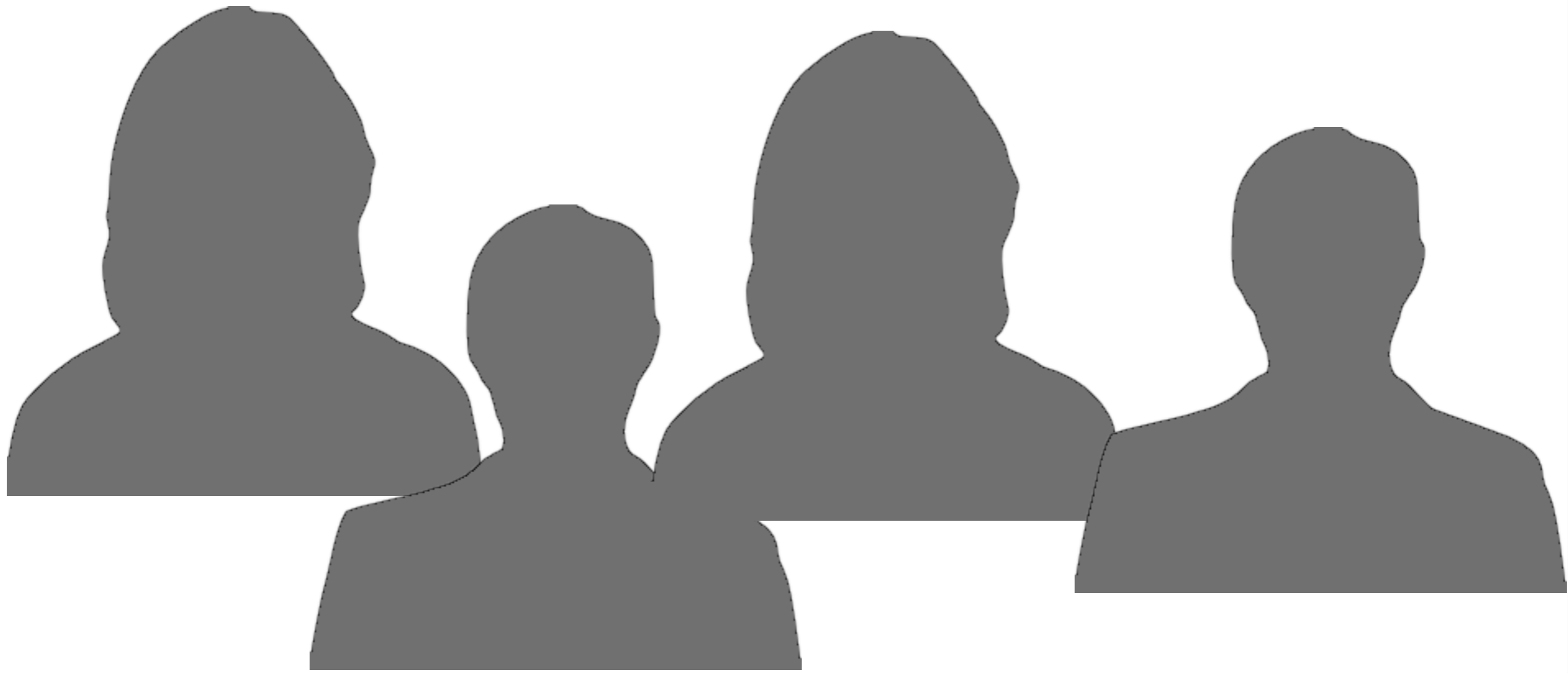
Trust

Economy

The Long Tail

Remixability

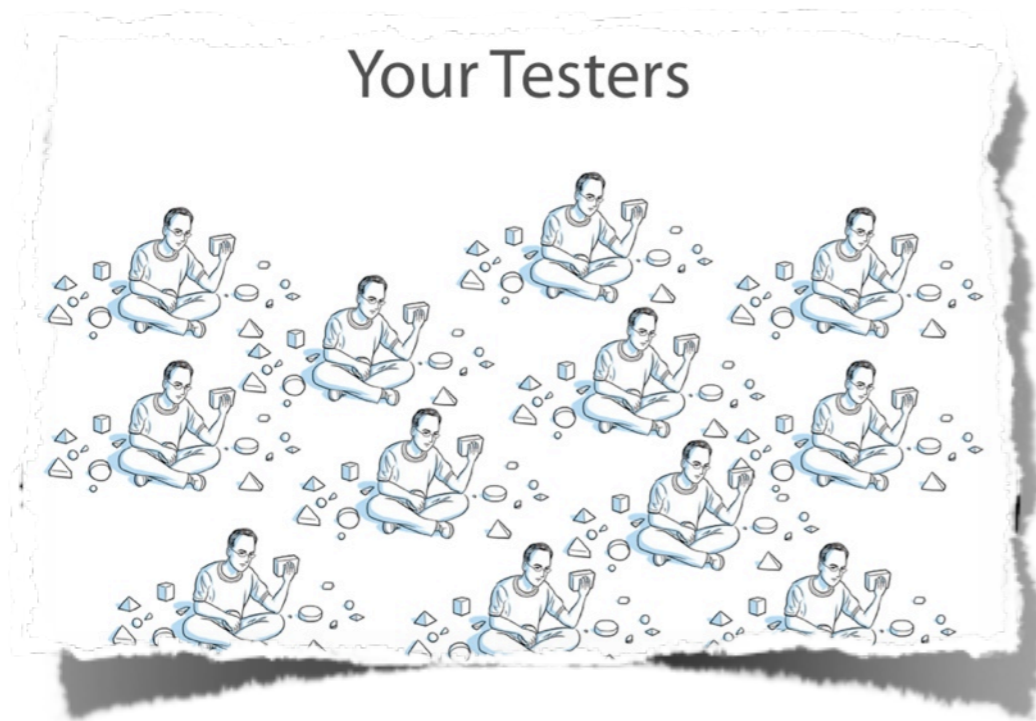
DataDriven



*Bachelor / Master Theses
in software mining*

Summary

Summary

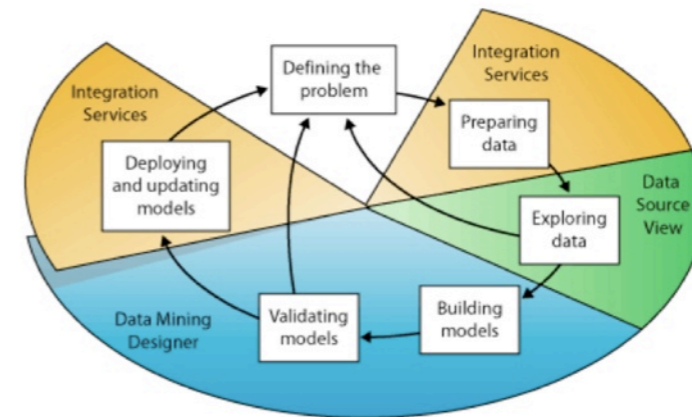


Summary

Your Testers



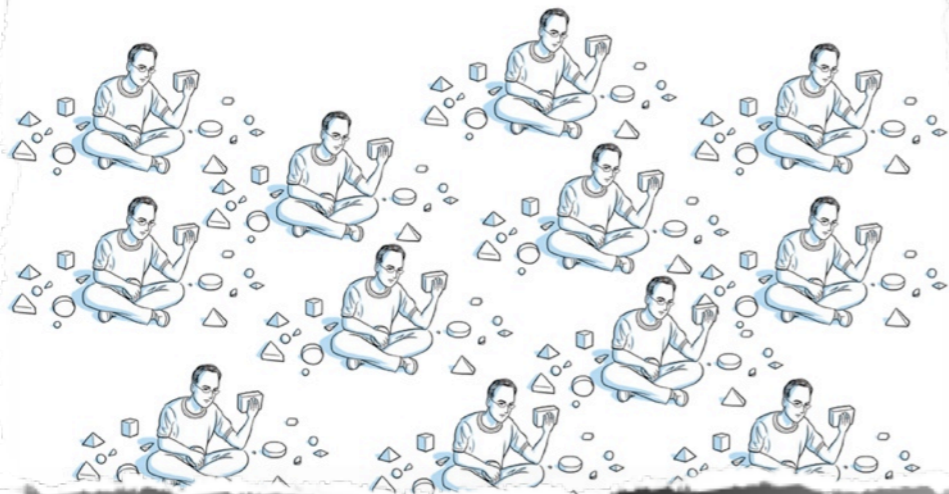
Building a mining model



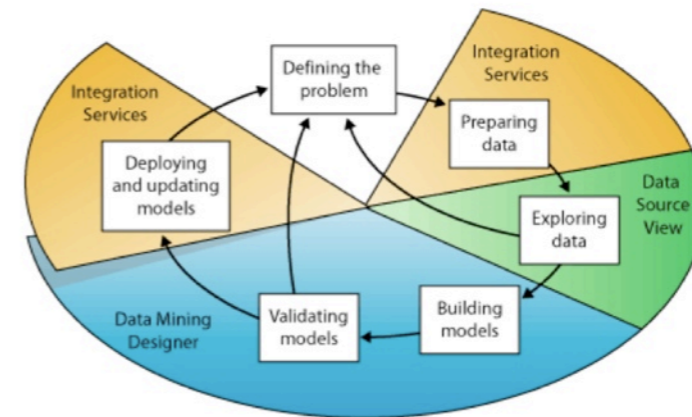
Source of image: <http://technet.microsoft.com/en-us/library/ms174949.aspx>

Summary

Your Testers



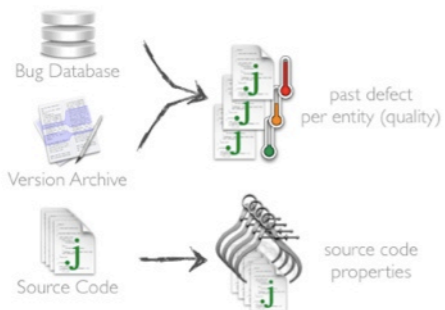
Building a mining model



Source of image: <http://technet.microsoft.com/en-us/library/ms174949.aspx>

Building the Model

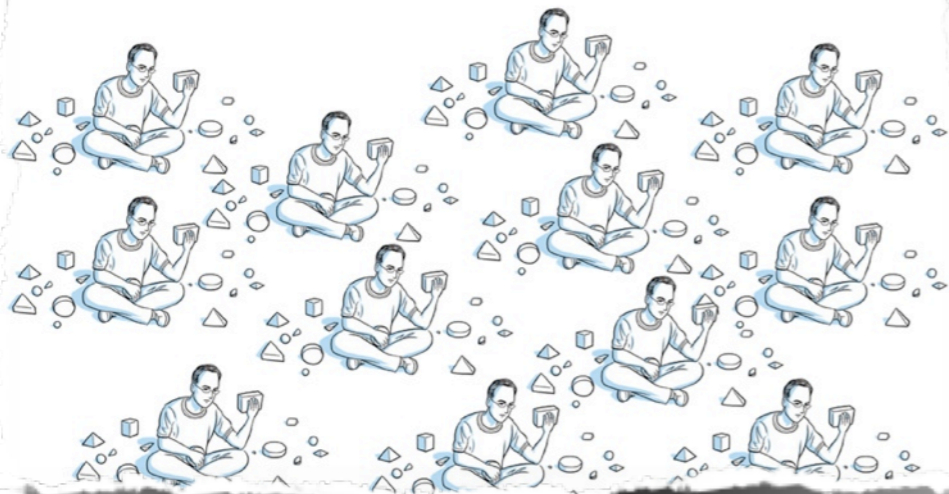
- Step 1: Define the problem
- Step 2: Prepare Data
- Step 3: Explore Data
- Step 4: Building the Model**
- Step 5: Validating the Model



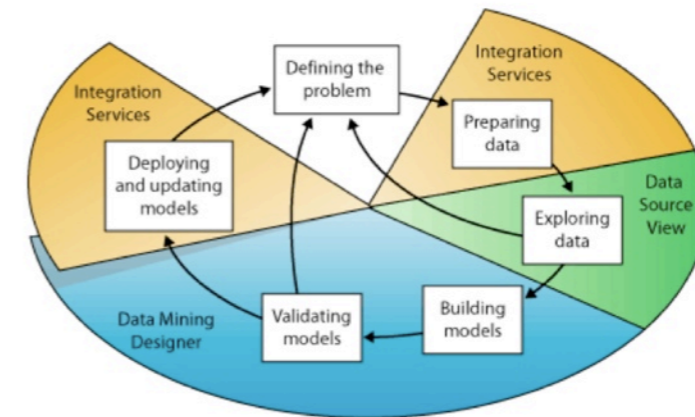
Do defects and metrics correlate?

Summary

Your Testers



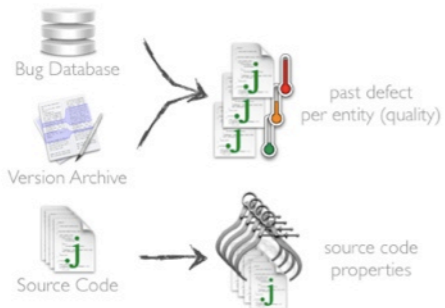
Building a mining model



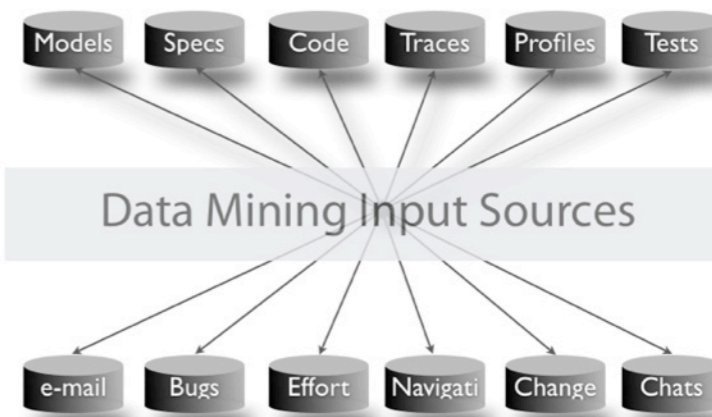
Source of image: <http://technet.microsoft.com/en-us/library/ms174949.aspx>

Building the Model

- Step 1: Define the problem
- Step 2: Prepare Data
- Step 3: Explore Data
- Step 4: Building the Model**
- Step 5: Validating the Model



Do defects and metrics correlate?

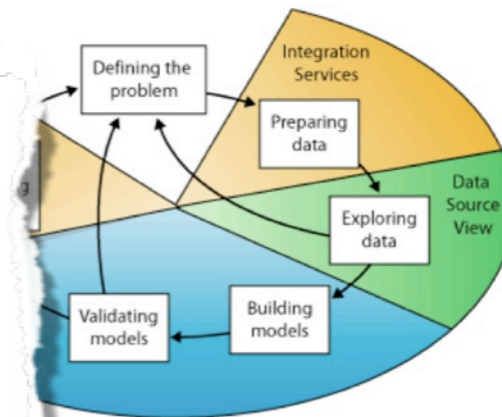


Summary

Your Testers



Building a mining model

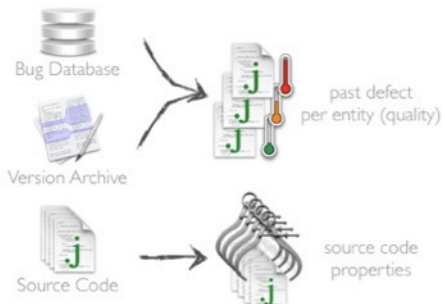


technet.microsoft.com/en-us/library/ms174949.aspx

Wikis
 Participation
 Recommendation
 Social Software
 Collaboration
 Perpetual Beta
 Joy of Use
 Usability
 Simplicity
Empirical SE 2.0
 Trust
 Economy
 The Long Tail
 Remixability
 DataDriven

Building the Model

- Step 1: Define the problem
- Step 2: Prepare Data
- Step 3: Explore Data
- Step 4: Building the Model**
- Step 5: Validating the Model



Do defects and metrics correlate?

