

More About Bound Variables

Substitution

Our definition of evaluation is based on the “substitution” of values for free variables within terms.

$$(\lambda x. t_{12}) v_2 \longrightarrow [x \mapsto v_2]t_{12} \quad (\text{E-APPABS})$$

But what is substitution, exactly? How do we define it?

Substitution

For example, what does

$$(\lambda x. x (\lambda y. x y)) (\lambda x. x y x)$$

reduce to?

Note that this example is not a “complete program” — the whole term is not closed. We are mostly interested in the reduction behavior of closed terms, but reduction of open terms is also important in some contexts:

- ▶ program optimization
- ▶ alternative reduction strategies such as “full beta-reduction”

Formalizing Substitution

Consider the following definition of substitution:

$$[x \mapsto s]x = s$$

$$[x \mapsto s]y = y$$

if $x \neq y$

$$[x \mapsto s](\lambda y. t_1) = \lambda y. ([x \mapsto s]t_1)$$

$$[x \mapsto s](t_1 t_2) = ([x \mapsto s]t_1)([x \mapsto s]t_2)$$

What is wrong with this definition?

Formalizing Substitution

Consider the following definition of substitution:

$$[x \mapsto s]x = s$$

$$[x \mapsto s]y = y$$

if $x \neq y$

$$[x \mapsto s](\lambda y. t_1) = \lambda y. ([x \mapsto s]t_1)$$

$$[x \mapsto s](t_1 t_2) = ([x \mapsto s]t_1)([x \mapsto s]t_2)$$

What is wrong with this definition?

It substitutes for free and *bound* variables!

$$[x \mapsto y](\lambda x. x) = \lambda x. y$$

This is not what we want!

Substitution, take two

$$[x \mapsto s]x = s$$

$$[x \mapsto s]y = y$$

if $x \neq y$

$$[x \mapsto s](\lambda y. t_1) = \lambda y. ([x \mapsto s]t_1)$$

if $x \neq y$

$$[x \mapsto s](\lambda x. t_1) = \lambda x. t_1$$

$$[x \mapsto s](t_1 t_2) = ([x \mapsto s]t_1)([x \mapsto s]t_2)$$

What is wrong with this definition?

Substitution, take two

$$[x \mapsto s]x = s$$

$$[x \mapsto s]y = y$$

if $x \neq y$

$$[x \mapsto s](\lambda y. t_1) = \lambda y. ([x \mapsto s]t_1)$$

if $x \neq y$

$$[x \mapsto s](\lambda x. t_1) = \lambda x. t_1$$

$$[x \mapsto s](t_1 t_2) = ([x \mapsto s]t_1)([x \mapsto s]t_2)$$

What is wrong with this definition?

It suffers from *variable capture*!

$$[x \mapsto y](\lambda y. x) = \lambda x. x$$

This is also not what we want.

Substitution, take three

$$[x \mapsto s]x = s$$

$$[x \mapsto s]y = y$$

if $x \neq y$

$$[x \mapsto s](\lambda y. t_1) = \lambda y. ([x \mapsto s]t_1)$$

if $x \neq y, y \notin FV(s)$

$$[x \mapsto s](\lambda x. t_1) = \lambda x. t_1$$

$$[x \mapsto s](t_1 t_2) = ([x \mapsto s]t_1)([x \mapsto s]t_2)$$

What is wrong with this definition?

Substitution, take three

$$[x \mapsto s]x = s$$

$$[x \mapsto s]y = y$$

if $x \neq y$

$$[x \mapsto s](\lambda y. t_1) = \lambda y. ([x \mapsto s]t_1)$$

if $x \neq y, y \notin FV(s)$

$$[x \mapsto s](\lambda x. t_1) = \lambda x. t_1$$

$$[x \mapsto s](t_1 t_2) = ([x \mapsto s]t_1)([x \mapsto s]t_2)$$

What is wrong with this definition?

Now substitution is a *partial function*!

E.g., $[x \mapsto y](\lambda y. x)$ is undefined.

But we want an result for every substitution.

Bound variable names shouldn't matter

It's annoying that that the “spelling” of bound variable names is causing trouble with our definition of substitution.

Intuition tells us that there shouldn't be a difference between the functions $\lambda x.x$ and $\lambda y.y$. Both of these functions do exactly the same thing.

Because they differ only in the names of their bound variables, we'd like to think that these *are* the same function.

We call such terms *alpha-equivalent*.

Alpha-equivalence classes

In fact, we can create equivalence classes of terms that differ only in the names of bound variables.

When working with the lambda calculus, it is convenient to think about these *equivalence classes*, instead of raw terms.

For example, when we write $\lambda x.x$ we mean not just this term, but the class of terms that includes $\lambda y.y$ and $\lambda z.z$.

We can now freely choose a different *representative* from a term's alpha-equivalence class, whenever we need to, to avoid getting stuck.

Substitution, for alpha-equivalence classes

Now consider substitution as an operation over *alpha-equivalence classes* of terms.

$$[x \mapsto s]x = s$$

$$[x \mapsto s]y = y$$

if $x \neq y$

$$[x \mapsto s](\lambda y. t_1) = \lambda y. ([x \mapsto s]t_1)$$

if $x \neq y, y \notin FV(s)$

$$[x \mapsto s](\lambda x. t_1) = \lambda x. t_1$$

$$[x \mapsto s](t_1 t_2) = ([x \mapsto s]t_1)([x \mapsto s]t_2)$$

Examples:

- ▶ $[x \mapsto y](\lambda y. x)$ must give the same result as $[x \mapsto y](\lambda z. x)$. We know the latter is $\lambda z. x$, so that is what we will use for the former.
- ▶ $[x \mapsto y](\lambda x. z)$ must give the same result as $[x \mapsto y](\lambda w. z)$. We know the latter is $\lambda w. z$ so that is what we use for the former.

Review

So what does

$(\lambda x. x (\lambda y. x y)) (\lambda x. x y x)$

reduce to?